

1 **Response to referee 1**

2 RC – Referee Comment

3 AC – Author Comment

4 MC – Manuscript change

5

6 **RC** - The paper untitled “An approach to the verification of high-resolution ocean models using spatial
7 methods” described a really interesting method to quantify benefit from high resolution model. The paper
8 describes in detail methodology and apply it to compare two ocean circulation forecast models on the
9 Nordic Sea. Scientific results obtain comparing the two forecast system are poorly commented and
10 explained but this scientific analysis is not the main topic of the paper, which is really dedicated to the
11 description, implementation of this methodology that was not already applied for ocean forecast. That
12 could be frustrating for readers, authors can certainly add analysis of some results, some suggestions are
13 provided below. Nevertheless, the paper is clear and objectives are well presented and I recommend the
14 publication of this paper if authors take into account few following remarks and comments.

15

16 **AC – We thank the reviewer for their time and expertise in reviewing the manuscript. Below are**
17 **our responses which we hope address the points raised along with changes made to the original**
18 **manuscript.**

19

20

21 **RC** - 1. Section 2, Figure 1 : this figure presents the domain and the difference of coastline between the
22 two models. Difference of coastline is an important point discussed also latter in the paper and illustrated
23 on fig 4. To be really interesting, I recommend to highlight the differences between the two SST fields on
24 this figure. A more contrasted color bar, for example, can highlight difference of spatial scale, intensity of
25 SST fronts. . . which are the main reasons to apply the HiRA method in this context.

26

27 **AC – Agreed. We looked at several different colour palettes which were also colour blind friendly**
28 **and replotted. In addition, some bathymetry contours were added to address a comment from**
29 **reviewer 2**

30 **MC – New colour scheme used, and bathymetry contours added.**

31

32

33

34 **RC** - 2. Section 3, line 187. Reference to WMO manual is useful but Authors should explained that this
35 guide refers to Atmosphere and that ocean scales are really different. In this paragraph specificities of
36 ocean should be described as difference of scale de-pending of the areas, open ocean vs shelf, rossby
37 radius . . . This is briefly discussed later in the section (line 245) but it should appear before in the
38 introduction of the method to justify to use it for ocean application.

39

40 **AC – Thank you, the WMO reference was indeed only atmospheric, tying in the original**
41 **justification and application of this method when it was applied to the atmosphere. As such we**
42 **have expanded the original section to refer to ocean specific characteristics, as well as a brief**
43 **addition to the introduction.**

44 **MC –** "A similar principle applies to the ocean, i.e. observations can represent an area around the
45 nominal observation location, though the representative scales are likely to be very different from
46 in the atmosphere. The representative scale for an observation will also depend on local
47 characteristics of the area, for example whether the observation is on the shelf, or in open ocean
48 or likely to be impacted by river discharge."

49

50 **RC -3.** Section 3, fig 3 and 4. Figure 3 and 4 are useful to understand the method and the neighbourhood
51 concept. But it could be really useful to have, on these figures or with a new figure, a clear description
52 (with an example) of how is computed the probability/density function especially in the coastal cases, how
53 the observations are selected in a neighbourhood, where the coastline is different between the two
54 models and when observations are removed from the statistics. A schematic view of this process should
55 be really useful to understand easily some non-intuitive results as for example why there is less
56 observation in a larger domain.

57

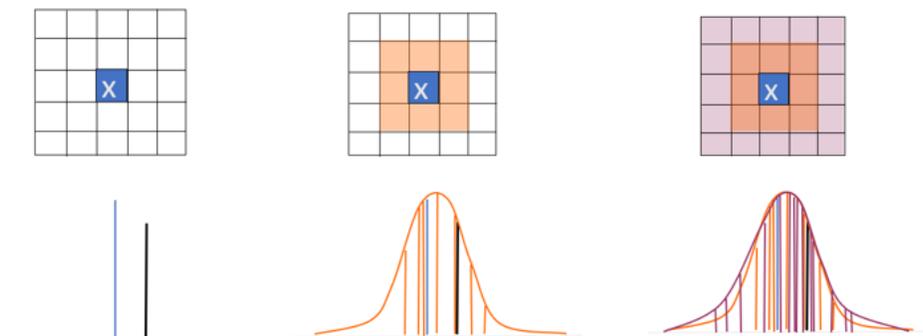
58 **AC –** We have added a schematic showing how the neighbourhood points contribute to
59 generating a pdf. We have also expanded the description of how missing points are handled
60 within the text.

61

62 **MC –** added as figure 4

63

64



65

66

67

68 **RC - 4.** Section 4, line 290. I suggest to use zonal and meridional instead of horizontal and vertical

69

70 **AC – Accepted**

71 **MC – Changed to zonal and meridional**

72

73

74 **RC - 5.** Section 4, figure 4. Unclear or a mistake in the legend. Why a) is 7x7 neighbourhood (NB4) and
75 b) NB5? Comparison should be done between similar neighbourhood.

76

77 **AC –** The idea we were trying to convey was that due to the forecast grids, the kilometre size of
78 neighbourhoods becomes increasingly incorrect as the neighbourhood becomes bigger if simply
79 assuming that multiplying 1.5 km or 7 km are accurate measures of the total size (instead of using

80 the true grid resolution in degrees). Coupled with that is the fact that the model resolution is
81 different in latitudinal and longitudinal directions.

82 **MC – We have separated out and expanded the table describing the neighbourhoods to indicate**
83 **why a 25x25 AMM15 is more suitable to match to 7x7 AMM7 than the 33x33 AMM15. Also modified**
84 **the caption to figure 4.**

85
86
87 **RC -6.** Section 5 and 6, fig 5, 7, 9, 10. It's really difficult to identify differences between each line,
88 probably too much lines on the same figures or more important line should be highlighted (in bold or with
89 darker color?) NB1 and NB2 are the more important, is really difficult to distinguished them especially on
90 fig7,9,10. Uncertainty, computed for each line, is difficult to associated to the right line. Is it useful to have
91 the "1" line for AMM15, there is no comparison with AMM7? It's also difficult on these figures to have clear
92 relationship between the uncertainty vertical bar and the difference bar. It will be useful to have on the
93 figure or in a table the information where the difference bars are smaller than the uncertainty. This is
94 discussed in the text (paragraph line 420) but it is difficult to verify what is described on the figures.

95
96 **AC – Agreed. We felt there was a balance to strike between showing how the scores change with**
97 **neighbourhood size and the ability to see detail of the actual results. The "1" is important in this**
98 **case as it shows the default result we would get if HiRA were not being used. However we have**
99 **tried to make the plots clearer whilst retaining that information.**

100 **MC - In order to clarify the plots we have removed some of the larger neighbourhoods from**
101 **figures 5, 9 and 10. In addition on figure 7 the main lines have been made bold. In order to help**
102 **with identifying where difference bars are less than the uncertainty, an S has been added over the**
103 **difference bars where the 95% uncertainly error bars of the two equivalent lines do not overlap.**

104
105 **RC -7.** Section 5. Discussion on the different results obtained on-shelf and off-shelf is really interesting,
106 but in the paper it appears as a mix between feasibility and useful methodology to compare several
107 forecasts and a clear difference due to dynamics, physical ocean process and seasonal cycle. I suggest
108 adding more quantitative information concerning the impact of the number of observation to compute
109 robust statistics. The sentence (line 460) explains that the model are better to forecast open ocean, but is
110 there any impact of the number of observation in the statistics? Do you compute statistics with the same
111 number of observations in the two domain (off-shelf and on-shelf)? Fig 12 and 13 seems to exhibit larger
112 uncertainty in the statistic on-shelf in comparison to off-shelf. On fig 12 and 13, it's clear that main
113 differences between the two models appears in summer. That's not really discussed in the paper, is there
114 clear explanation, is it due to physical seasonal processes or mainly due to the number of observations?

115
116 **AC – The aim of this paper was to show how the HiRA technique could be used to tease out**
117 **interesting detail of the model forecasts which could then be a basis for investigation in the**
118 **future. Notably in this case the apparent seasonal signal. Figure 14 indicates the numbers of**
119 **observations going into the two domains, and hence the fact that this is a potential source of**
120 **error. However with the underlying characteristics of the domains being different, it is quite likely**
121 **that the spatial distribution of observations within the domains is as important as the number of**
122 **observations. Again, as this was meant to be an investigation of the potential for the verification**
123 **technique rather than a full model assessment, we did not dig further into the detail of this, but do**
124 **think it is an important consideration when assessing any results produced using HiRA**

125
126 **RC -8.** Section 5, line 479. Conclusion of this paragraph is not clear. What do you really mean by "closer
127 look at the data"?

128 **AC – Essentially breaking the data down and identifying underlying specific parts of the data**
129 **which may be contributing to the results counter to the general trend, and which are masked by**
130 **aggregating.**

131 **MC – Edited the text to “therefore a deeper look at the data is required to assess whether this**
132 **signal is consistent within shorter time periods, or whether there are underlying periods**
133 **contributing significant and contrasting results to the whole-period aggregate. “**

134

135 **RC -9.** Section 5, line 508. Last sentence concludes on differences at NB2 scale, could you add comment
136 on this conclusion about significance and robustness of this result.

137 **AC – Yes, as you indicate the statement is too strong given the error bars presented. We have**
138 **highlighted that aspect (indicating that the error bars cross, so whilst we cannot say that the**
139 **difference is significant, we cannot, with the plot provided, say they are not.) And giving**
140 **suggestions as how to improve this.**

141 **MC – “At the NB2 scale, the AMM15 potentially demonstrates more benefit than AMM7 except for**
142 **April and May, where the two show similar results. There is a balance to be struck in this**
143 **conclusion as the differences between the two models are rarely greater than the 95% error bars.**
144 **This in itself does not mean that the results are not significant. However, care should be taken**
145 **when interpreting such a result as a statistical conclusion rather than broad guidance as to model**
146 **performance. Attempts to reduce the error bar size, such as increasing the number of**
147 **observations, or number of times within the period would aid this interpretation.”**

148

149 **RC -10.** Section 5, fig 14. On this figure lack of observations seems to appear end of May and in the text
150 (line 487) authors indicate that missing data are in April.

151 **AC – The text was incorrect, there was a reduction in observations during May due to issues with**
152 **the observation extraction from CMEMS. Additionally, there was a forecast reduction in April (due**
153 **to separate technical issues) not indicated by the plot.**

154 **MC – The text is now correct and additionally refers to the missing forecast period.**

155

156

157

158

159

160

161

162

163

164

165 Response to Referee 2

166 RC – Referee Comment

167 AC – Author Comment

168 MC – Manuscript change

169

170 **RC** - In this contribution, the authors conduct a skill assessment of two operational ocean
171 models running in the North West European Shelf with different configurations and spatial
172 resolutions. Since the increased spatial resolution might require ad hoc metrics to properly
173 reflect the model performance and reduce the impact of so-called double-penalty effects
174 (occurring when using point-to-point comparisons with features present in the model but
175 misplaced with respect to the observations), the present work is welcomed. It addresses this
176 interesting and essential topic by intercomparing models' performances in overlapping regions
177 to infer their respective strengths and weaknesses. Equally, the methodology proposed is
178 consistent and the results obtained are relevant, especially in the framework of the Copernicus
179 Marine Service (albeit not explicitly mentioned in the document)

180 Based on my expertise on ocean models validation, I particularly appreciate the proposed
181 approach (named HiRA) since it might be useful in parent-son inter-comparisons in order to
182 quantify the added-value of downstream services such as very high resolution coastal models
183 (embedded into CMEMS regional ocean forecasting systems) that are currently running in port-
184 approach areas. I am confident this work can attract the interest of the scientific audience,
185 being cited in future works dealing with similar issues. The style was fluent although some parts
186 (mainly the introduction and the references) could be revised and enhanced. Based on my
187 judgment, I deem the manuscript acceptable upon minor revision. In the following lines I
188 provide some comments, which should hopefully strengthen even more the manuscript.

189 **AR – We thank reviewer 2 for their time and effort reviewing this paper, which produced**
190 **some very interesting comments. Below we have addressed each point and hope this has**
191 **resulted in a stronger paper.**

192

193

194

195 **RC** - General comments:

196 1. Since the main purpose of this work is to showcase the potential of the proposed
197 methodology in operational ocean forecasting, I miss a reference to the Copernicus Marine
198 Environment Monitoring Service -CMEMS- (Le Traon et al., 2019)., although the in situ

199 observations used here were downloaded from CMEMS catalogue. Within this context, there
200 are some valuable and concerted initiatives such as the Product Quality Working Group (PQWG)
201 or the North Atlantic Regional VALidation tool -NARVAL-(Lorente et al., 2019) where physical
202 and biogeochemical model intercomparisons are conducted on a regular basis to deliver
203 outcomes to a broad scientific community.

204 Le-Traon et al.,2019.“From Observation to Information and Users:The Copernicus Marine
205 Service Perspective”.Front.Mar.Sci., 22,https://doi.org/10.3389/fmars.2019.00234.

206 Lorente et al., 2019. “The NARVAL software toolbox in support of ocean models skill
207 assessment at regional and coastal scales”. Computational Science, ICCS 2019.

208 **AR -Thankyou, these references have been added.**

209 **MC – Additional references have been included within the introduction of the paper.**

210

211 **RC - 2.** Equally, I also miss a reference to GODAE Coastal Ocean and Shelf Seas Task Team
212 (COSS-TT), where the Met Office is an active member, involved in a wealth of valuable
213 initiatives in terms of ocean model inter-comparisons. In this context, I think that the state-of-
214 art about previous inter-comparison exercises is not thorough and is poorly cited, despite of the
215 abundant literature reported elsewhere. In this work, there are only 28 references (which is
216 insufficient) and nearly the 50% of them were published in 2010 or earlier, so an update is
217 highly recommended. Below I suggest a number of recent works to build upon:

218 Aznar et al, 2015. “Strengths and weaknesses of the CMEMS forecasted and reanalyzed
219 solutions for the Iberia-Biscay-Ireland (IBI) waters”. Journal of Marine Systems,159, 1-14.

220 Mourre et al., 2019. “Assessment of High-Resolution Regional Ocean Prediction Systems Using
221 Multi-Platform Observations: Illustrations in the Western Mediterranean Sea”.

222 Lorente et al., 2019. “Skill assessment of global, regional, and coastal circulation fore-cast
223 models: evaluating the benefits of dynamical downscaling in IBI (Iberia–Biscay–Ireland) surface
224 waters”. Ocean Science, 15, 967-996. Doi: /10.5194/os-15-967-2019.

225 Mason et al., 2019. “New insight into 3-D mesoscale eddy properties from CMEMS operational
226 models in the western Mediterranean”. Ocean Science, 15, 1111–1131.

227 Hernández et al., 2018. “Measuring performances, skill and accuracy in operational
228 oceanography: New challenges and approaches”. In "New Frontiers in Operational
229 Oceanography", Eds. GODAE OceanView, 759-796, doi:10.17125/gov2018.ch29.

230 Juza et al, 2015. “From basin to sub-basin scale assessment and intercomparison of numerical
231 simulations in the western Mediterranean Sea”. Journal of Marine System,149, 36-49,
232 doi:10.1016/j.jmarsys.2015.04.010.

233 Hernández et al., 2015. “Recent progress in performance evaluations and near real-time
234 assessment of operational ocean products”. Journal of Operational Oceanography, 8, Issue
235 sup2: GODAE OceanView Part 2

236 Rockel et al., 2015. “The regional downscaling approach: a brief history and recent advances”.
237 Curr. Clim. Change Rep., 1, 22–29, <https://doi.org/10.1007/s40641-014-0001-3>.

238 Katavouta et al, 2016. “Downscaling ocean conditions with application to the Gulf of Maine,
239 Scotian shelf and adjacent deep ocean”. Ocean Model., 104, 54–72.

240 And some other older works:

241 Crosnier, L., and C. Le Provost. 2007. “Inter-comparing five forecast operational systems in the
242 North Atlantic and Mediterranean basins: The MERSEA-strand1 method-ology”. Journal of
243 Marine Systems, 65, 354–375.

244 Greenberg et al, 2007. “Resolution issues in numerical models of oceanic and
245 coastal circulation”. Cont. Shelf Res., 27, 1317–1343.

246 Hernández, 2011. “Performance of Ocean Forecasting Systems—Intercomparison Projects” u.
247 Book: Operational Oceanography in the 21st Century, Chapter 23.

248 **AR – Thank you for the additional references. A selection of these have been added to the**
249 **text in the introduction section to broaden the description of the existing state of things.**

250 **MC – Additional references have been included within the introduction of the paper.**

251

252

253 **RC - 3.** In section 1 (Introduction), a preliminary paragraph about why model inter-comparisons
254 are necessary would be convenient. Equally, a brief description of the types of inter-
255 comparisons exercises would be pertinent:

256 i) between two different forecasting systems in the overlapping region to check the consistency
257 of each model solution;

258 ii) between two versions of the same system, in order to evaluate the added-value of the
259 upgraded one before it is transitioned into fully operational status;

260 iii) a parent-son inter-comparison, to evaluate the quality of the downscaling approach
261 adopted;

262 iv) a comparison between both the forecast and the reanalyzed solutions of the same model in
263 order to infer the primary role of both the grid resolution and the atmospheric forcing,
264 especially in coastal areas (see Aznar et al., 2015, for further details).

265 **AR - We have introduced this in combination with some of the references in RC2.**

266

267 **RC - 4.** In section 2.1 (Data and Methods: Forecast), I strongly suggest adding a table to provide
268 a general overview of the two model's main features in a more synthesized way: version of
269 model, geographic domain, grid resolution, number of depth levels, number of forecast
270 horizons, open boundary conditions, tidal forcing, atmospheric forcing, river forcing,
271 assimilation scheme, bathymetry, etc. Although most of this information is already provided in
272 the text, I think a table would be rather useful as a summary.

273 **AR - We have added a summary table of the differences relevant for this study.**

274

275 **MC - Additional table (table 1) added in the manuscript**

276

277 **RC- 5.** In section 2.1 (Data and Methods: Forecast), neither river forcing is mentioned, nor river
278 freshwater discharge is taken into account when describing the general considerations. The
279 study-area comprises several rivers estuaries (Seine, Rhine, even Loire) with significant
280 freshwater runoff that might eventually impact on the SST field in coastal areas. Figure 2 shows
281 that some stations are located quite close to those rivers mouth. Please clarify this point, why
282 the river forcing is out of the discussion. In particular, Graham et al (2018) suggested that
283 AMM7 configuration might be more diffusive than AMM15 within river plumes, allowing
284 freshwater input from the Rhine to be advected offshore.

285 **AR - Even if it is true that the river forcing plays a role in the coastal areas, it has been proven**
286 **in Tonani et al. 2019 that it has a very small impact on SST. It is much more evident in surface**
287 **salinity. We describe in 2.1 the characteristics/differences that are relevant for this study, a**
288 **comprehensive description of the two forecasting systems is in Tonani et al. 2019. We specify**
289 **in section 5 that this study is not focused on the coastal areas due to the assumptions in the**
290 **choice for the neighbour, with different number of observations in the two configurations**
291 **due mainly to the Land-Sea mask differences.**

292 **This is a very interesting issue and there is a need for a coastal-focused assessment of the**
293 **forecast. We will take it account for future work. It is also worth to notice that comparing the**
294 **model only simulation (non-assimilative analysis) over a long period (30 years) of Graham et**
295 **al. 2018 with 9 months of assimilative analysis-forecast could be misleading for the reader.**
296 **Graham et al. experiment is using different lateral boundaries and a different atmospheric**
297 **forcing and doesn't have data assimilation. Both AMM7 and AMM15 forecasting systems are**
298 **assimilating SST obs (Insitu and satellite). Even if we consider negligible these differences, it's**
299 **difficult to justify the comparison of a seasonal mean over 30 years with few months of**
300 **forecast. The results from Graham et al. 2018 have not been confirmed by Tonani et al. 2019,**

301 while assessing the operational trials (with data assimilation and the operational forcing as
302 described in the paper) against OSTIA. This validation is shown at basin level in the paper, but
303 we did analyse also off-shelf and on-shelf differences. There are no significant differences
304 compared to the full domain inter-comparison.

305 From Tonani et al. 2019:

306 *“Temperature RMSD and bias are very small at surface, due to the strong constraint of the data*
307 *assimilation of SST (as described in 4.3) while at the bottom AMM15 is more accurate in prescribing*
308 *the temperature at all mooring locations (Error! Reference source not found.).*

309 *AMM7 and AMM15 both have high salinity errors in the German Bight, as highlighted by the*
310 *comparison with the buoys that are located closer to the coast (Fino1, Fino3 and UFSDeBucht). This is*
311 *most probably due to representation of river discharge. AMM15 performs better than AMM7,*
312 *probably because it is less diffusive within river plumes and has a lower lateral diffusion. Improved*
313 *bathymetry and coastal resolution are also likely to play a role in coastal areas with depth less than*
314 *20m. AMM15 has halved the salinity error compared to AMM7 when compared with the outer buoys*
315 *(NsbII and TWEms). It is encouraging to see that AMM15 is better than AMM7 at the bottom at all*
316 *mooring locations. The decision to use the climatological river discharge dataset instead of E-Hype for*
317 *AMM7, and subsequently AMM15, has improved salinity remarkably in the German Bight, reducing*
318 *the model fresh bias. This modification was implemented in April 2017, meaning that we have*
319 *significantly improved the salinity in the last two major updates of the NWS forecasting system.*
320 *Nevertheless, using a climatological river runoff dataset is a limitation for a high-resolution*
321 *forecasting system, affecting variability in coastal water properties. Finding a suitable alternative will*
322 *be a priority for future releases of this system.”*

323 **MC - No changes**

324

325 **RC- 6.** In the same line, an event-oriented inter-comparison (with a focus on river plumes and
326 abrupt SST drops due to impulsive-type riverine discharges) would allow you to better infer the
327 ability of each system to capture small-scale coastal processes (with and without HiRA
328 approach). This process-based validation approach, albeit commonly used in meteorology and
329 weather forecasting, is rather novel in operational oceanography and mostly devoted to
330 extreme sea level and wave height episodes. I am not asking to provide new and
331 complementary analysis but please take it as a kind suggestion for future works.

332 **AR - Yes, we agree. We will take this comment into consideration for future work.**

333 **MC - no changes**

334

335 **RC- 7.** With regards to the double-penalty effect, I was somehow expecting a multi-parameter
336 analysis, with a special focus on altimetry products, sea level anomalies and mesoscale eddies.
337 Did you have the chance to test HiRA approach with other variables? If so, could you add a

338 comment about it, even if you only obtained preliminary results? If not, I think this task should
339 remain as a priority for future works and thus be explicitly mentioned in the text.

340 **AR - Within this assessment we started simple, since we wanted to know whether the**
341 **technique had anything to offer and only looked at SST, though other parameters were**
342 **considered (e.g. velocities). One of the next steps will be to apply this to a broader range of**
343 **parameters. – We have noted this in the conclusions.**

344 **MC – The conclusions section has been updated.**

345 *The conclusion does mention that other variables should be assessed. (lines ~631) – Agree that*
346 *more parameters would be good.*

347

348 **RC- 8.** Likewise, I miss a deeper discussion respect to the previous works by Tonani et al(2019)
349 and especially that one by Graham et al (2018) where a “traditional point-to-point SST
350 validation approach” was performed with the new AMM15 system. I think that the fact of
351 contrasting results from both papers / both methodologies could benefit the discussion section,
352 particularly when dealing with on-shelf and off-shelf differences as far as Graham et al (2018)
353 proved the reduction in seasonal SST bias was greater off-shelf than on-shelf when using
354 AMM15 (which supports the results exposed in Figures 9 and 10 of the present work). Again,
355 on-shelf results were worse and you succinctly listed river mixing as a potential source of
356 uncertainties, but no additional information was provided about the role of river forcing (as I
357 aforementioned in point5). I guess that the river fluxes could have been altered between the
358 two models (being one configuration fresher and cooler than the other).

359 **AR – Please see also answer to comment 5. The work of this paper is focused on 9 months of**
360 **forecast validation, while the study discussed in Graham et al 2018 is based on a model only**
361 **simulation over 30 years. The SST data assimilation has significant impact on both systems**
362 **(AMM7 and AMM15) due to the good coverage of the observations. The comparison**
363 **between results of a model only long simulation against few months of forecasts is not**
364 **straightforward and implies several assumptions that deviates from the object of this paper**
365 **to assess the forecasts skills in different configurations. We explained in answer 5 that the**
366 **rivers seem to play a minor role on SST and that we need a specific study focused on the**
367 **coastal area. The freshwater inflow has for sure an important impact on the stratification and**
368 **this needs to be properly assessed. The differences on the freshwater are due to horizontal**
369 **and vertical resolution. Bathymetry and model diffusivity. It is a complex combination of**
370 **different aspects that is not addressed in this work.**

371 **MC – no change**

372

373 Minor comments:

374 **RC** - Abstract: I recommend explaining briefly (in two lines) the double penalty effect as part of
375 the potential audience might not be familiarized with this concept. For instance: “[...]referred
376 to as the double-penalty effect, occurring in point-to-point comparisons with features present
377 in the model but misplaced with respect to the observations.”

378 **AR – Added brief explanation**

379 **MC** - “...the double-penalty effect. This effect occurs in point-to point comparisons whereby
380 features correctly forecast but misplaced with respect to the observations are penalised
381 twice; once for not occurring at the observed location, and secondly for occurring at the
382 forecast location, where they have not been observed.”

383

384 **RC** – Keywords: I suggest adding “skill assessment”, “validation” and/or “double-penalty”.

385 **AR – Additional keywords will be added**

386 **MC – Added ‘double penalty’ and validation to keywords**

387

388 **RC** - Figure 1: As previously indicated by the anonymous reviewer 1, a more contrasted color
389 bar is required to highlight the spatial SST differences. Bathymetric contours would be also
390 welcomed.

391 **AR** - Agreed. We looked at a number of different colour palettes which were also colour blind
392 friendly and replotted.

393 **MC** – New colour scheme used and bathymetry contours added.

394

395 **RC** - Figure 8: Albeit rather obvious, please indicate that masked regions are in grey color.

396 **AR** - Agreed

397 **MC** - Added “data within the grey areas is masked” to caption

398

399 **RC** - Introduction: Lines 58-60: That sentence sounds odd. Could you rephrase it, please?

400 **AR** - Agreed.

401 **MC** – Added “In these methods forecasts are assessed at multiple spatial or temporal scales
402 to see how model skill changes as the scale is varied.”

403

404 **RC** - Line 61: please replace “suggested” by “suggesting”

405 **AR - Done**

406 **MC – “suggested” replaced by “suggesting”**

407

408 **RC - Line 65: please replace “more like” by “more similar to”**

409 **AR - Accepted**

410 **MC – replaced “more like” with “more similar to”**

411

412 **RC - Section 2.1. Forecast**

413 Lines 106-108: I guess that hourly instantaneous values are provided for the sea surface and

414 daily averages for the rest of the water column. Please, could you clarify it?

415 **AR – This is now clarified in the text**

416 **MC - “Hourly instantaneous values and daily 25-hour, de-tided, averages are provided for the**

417 **full water column. “**

418

419 **RC - Line 117: Why the study period comprises from January to September 2019? Any chance to**

420 **expand the analysis to cover the entire 2019 year? That would be interesting to infer seasonal**

421 **differences between both model configurations... –**

422 **AR – The study period could be expanded to cover a longer period, however since this was an**

423 **introductory study, we felt that the full benefit of a longer assessment period should also**

424 **involve additional parameters and a more focussed assessment of the model, rather than, as**

425 **here, an assessment of the method. The potential seasonal signal gives a focus to any further**

426 **study.**

427 **MC – no change**

428

429 **RC-Lines 132-133: please comment that semi-diurnal M2 is one of the predominant tidal**

430 **constituents in this region (that is the reason to compute means over 25 hours in order to**

431 **remove the tidal signal).**

432 **AR – The major tidal constituent over the North West European shelf is the semidiurnal lunar**

433 **component, M2. It has a period of 12 h 25 min (Howarth, M. and Pugh, D.: Chapter 4**

434 **Observations of Tides Over the Continental Shelf of North-West Europe, Elsevier**

435 **Oceanography Series, 35, 135–188, [https://doi.org/10.1016/S04229894\(08\)70502-6](https://doi.org/10.1016/S04229894(08)70502-6), 1983.).**
436 **The 25 hours mean is therefore removing (or filtering out) the tidal signal.**

437 **MC - “The tidal signal is removed because the period of the major tidal constituent, the**
438 **semidiurnal lunar component M2, is 12hr and 25min (Howarth and Pugh, 1983).**

439

440 **RC - Section 7: Discussion and conclusions.**

441 Lines 538-539: as previously indicated, provide further insight into on-shelf and off-shelf
442 differences, contrasting the results obtained with those reported in Graham et al (2018).

443 **AR - See also answers to comment 5 and 8.**

444 **MC - Added in the conclusion “Forecast verification studies tailored for the coastal/shelf**
445 **areas are needed for properly understand the forecast skills in areas with high complexity**
446 **and fast evolving dynamics.”**

447

448 **RC - Lines 540-545: is there any adopted rule or any agreed proposal to wisely select the**
449 **neighborhood sizes?**

450 **AR – There is no accepted rule for doing this, particularly as the appropriate neighbourhood**
451 **sizes will likely be different for different parameters. As the neighbourhoods become larger**
452 **the CRPS will (generally) become smaller, but the improvements in the score will occur in**
453 **smaller increments as the neighbourhood grows. However, there then comes a point where**
454 **the neighbourhood becomes too large and points are introduced into the neighbourhood for**
455 **which the observation (which is fixed at a point) is no longer representative, at which point**
456 **the CRPS tends to increase again (degrade). It is therefore possible to infer something about**
457 **the representativeness of the observation and the level of variability within the**
458 **neighbourhood. In a homogeneous field you could infer aspects of representativeness, but**
459 **because the sampling is rarely that homogenous (except perhaps deep ocean) it’s not easy to**
460 **infer anything general which can be applied all the time. We investigated whether there were**
461 **any specific scales that could be identified but could not definitively draw any conclusions.**
462 **The current method, when applied to the atmosphere, is to initially use a broad set of**
463 **neighbourhoods and then use a subset of these for routine verification once / if the**
464 **representativeness for a variable becomes apparent (i.e. the CRPS starts to degrade / tail off).**

465

466

467 An approach to the verification of high-
468 resolution ocean models using spatial methods

469 Ric Crocker¹, Jan [Maksymczuk](#)²~~Maksymczuk~~[±], Marion Mittermaier¹, Marina [Tonani](#)²~~Tonani~~[±],
470 Christine [Pequignet](#)²~~Pequignet~~[±]

471 ¹[Verification, Impacts and Post-Processing, Weather Science, Met](#)[±]~~Met~~ Office, Exeter, [EX1 3PB](#)~~EX1-3PB~~, [United](#)
472 [Kingdom](#)~~UK~~

473 ²[Ocean Forecasting Research & Development, Weather Science, Met Office, Exeter, EX1 3PB, United Kingdom](#)

474

475

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

477 Abstract

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

481 Whilst qualitative assessments have demonstrated the benefits brought by the increased
482 resolution of AMM15, particularly in the ability to resolve finer-scale features, it has been difficult
483 to show this quantitatively, especially in forecast mode. Application of typical assessment metrics
484 such as the root mean square error have been inconclusive, as the high-resolution model tends
485 to be penalised more severely, referred to as the double-penalty effect. [This effect occurs in
486 point-to point comparisons whereby features correctly forecast but misplaced with respect to
487 the observations are penalised twice; once for not occurring at the observed location, and
488 secondly for occurring at the forecast location, where they have not been observed.](#)

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

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

502 Keywords

503 verification, ocean forecasts, SST, spatial methods, neighbourhood, [validation, double-penalty](#)

504

505 1. Introduction

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

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

- 512 • between different forecasting systems over an overlapping region to check for model
513 consistency between the two;
- 514 • between two versions of the same model to test the value of model upgrades prior to
515 operational implementation;
- 516 • parent-son intercomparison, evaluating the impact of downscaling or nesting of models;
517 • a forecast comparison against reanalysis of the same model, inferring the effect of
518 resolution and forcing, especially in coastal areas.

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

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

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

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

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

546 Spatial methods allow forecast models to be assessed with respect to several different types of
547 focus. Initially these methods were classified into four groups. Some methods look at the ability
548 to forecast specific features (e.g. Davis et al., 2006), some look at how well the model performs
549 at different scales (scale-separation, e.g. Casati et al., 2004). Others look at field deformation
550 (how much a field would have to be transformed to match a ‘truth’ field (e.g. Keil and Craig,
551 2007). Finally, there is neighbourhood verification, many of which are equivalent to low band-
552 pass filters. ~~In these methods, whereby values of forecasts in spatio-temporal neighbourhoods~~
553 ~~are assessed to see at multiple what spatial or temporal scales to see how model scale certain~~
554 ~~levels of skill changes as the scale is varied are reached by a model.~~

555 Dorninger et al. (2018) provides an updated classification of spatial methods,
556 ~~suggesting suggested~~ a fifth class of methods, known as distance metrics, which sit between field
557 deformation and feature-based methods. These methods evaluate the distances between
558 features, but instead of just calculating the difference in object centroids (which is typical), the
559 distances between all grid point pairs are calculated, which makes distance metrics more [similar](#)

560 ~~to like~~ field deformation approaches. Furthermore, there is no prior identification of features.
561 This makes distance metrics a distinct group that warrants being treated as such in terms of
562 classification. Not all methods are easy to classify. An example of this is the Integrated Ice Edge
563 Error (IIEE) developed for assessing the sea ice extent (Goessling et al., 2016).

564 This paper exploits the use of one such spatial technique for the verification of sea surface
565 temperature (SST), in order to determine the levels of forecast accuracy and skill across a range
566 of model resolutions. The High-Resolution Assessment framework (Mittermaier, 2014,
567 Mittermaier and Csimá, 2017) is applied to the Met Office Atlantic Margin Model running at 7 km
568 (O’Dea et al., 2012, O’Dea et al., 2017, King et al., 2018) (AMM7), and 1.5 km (Graham et al.,
569 2018, Tonani et al., 2019) (AMM15) resolutions for the European North West Shelf (NWS). The
570 aim is to deliver an improved understanding beyond the use of basic biases and RMS errors for
571 assessing higher resolution ocean models, which would then better inform users on the quality
572 of regional forecast products. Atmospheric science has been using high-resolution convective-
573 scale models for over a decade, and so have experience in assessing forecast skill on these scales,
574 so it is appropriate to trial these methods on eddy-resolving ocean model data. As part of the
575 demonstration, the paper also looks at how the method should be applied to different ocean
576 areas, where variation at different scales occurs due to underlying driving processes.

577 ~~This paper will demonstrate one of these spatial frameworks, HiRA (Mittermaier, 2014), and~~
578 ~~apply it to sea surface temperature (SST) daily mean forecasts from the Met Office operational~~
579 ~~ocean systems for the European North West Shelf (NWS). As part of the demonstration, the paper~~
580 ~~also looks at how the method should be applied to different ocean areas, where variation at~~
581 ~~different scales occurs due to underlying driving processes.~~

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

586 Section 2 describes the model and observations used in this study along with the method applied.
587 Section 3 presents the results, and section 4 discusses the lessons learnt while using HiRA on

588 ocean forecasts and sets the path for future work by detailing the potential and limitations of the
589 method.

590

591 2. Data and Methods

592 2.1 Forecasts

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

- 596 • NORTHWESTSHELF_ANALYSIS_FORECAST_PHYS_004_001_b (AMM7)
- 597 • NORTHWESTSHELF_ANALYSIS_FORECAST_PHY_004_013 (AMM15)

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

609

	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

610

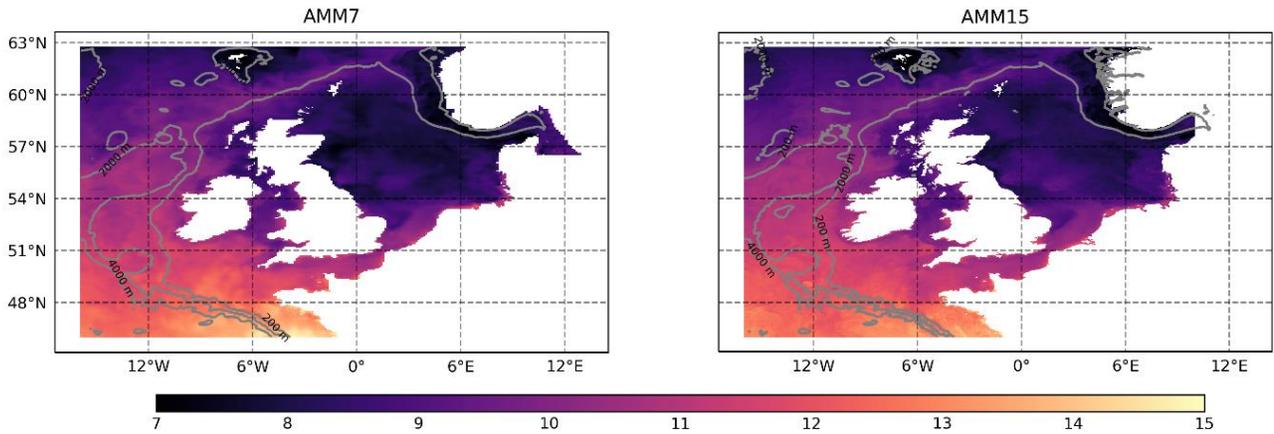
611 [Table 1: Summary of the main differences between NORTHWESTSHELF_ANALYSIS_FORECAST_PHYS_004_001_b \(AMM7\) and](#)
612 [NORTHWESTSHELF_ANALYSIS_FORECAST_PHYS_004_013 \(AMM15\)](#)

613 The AMM15 and AMM7 systems run once a day and provide forecasts for temperature, salinity,
614 horizontal currents, sea level, mixed layer depth, and bottom temperature. [Hourly](#)~~These~~
615 ~~products are provided as hourly~~ instantaneous [values](#) and daily 25-hour, de-tided, averages [are](#)
616 [provided for the full water column](#).

617 AMM7 has a regular latitude-longitude grid, whilst AMM15 is computed on a rotated grid and re-
618 gridded to have both models delivered to the (CMEMS) data catalogue
619 (<http://marine.copernicus.eu/services-portfolio/access-to-products/>) on a regular grid. A fuller
620 description of the respective configurations of the two models can be found in Tonani et al.,
621 (2019).

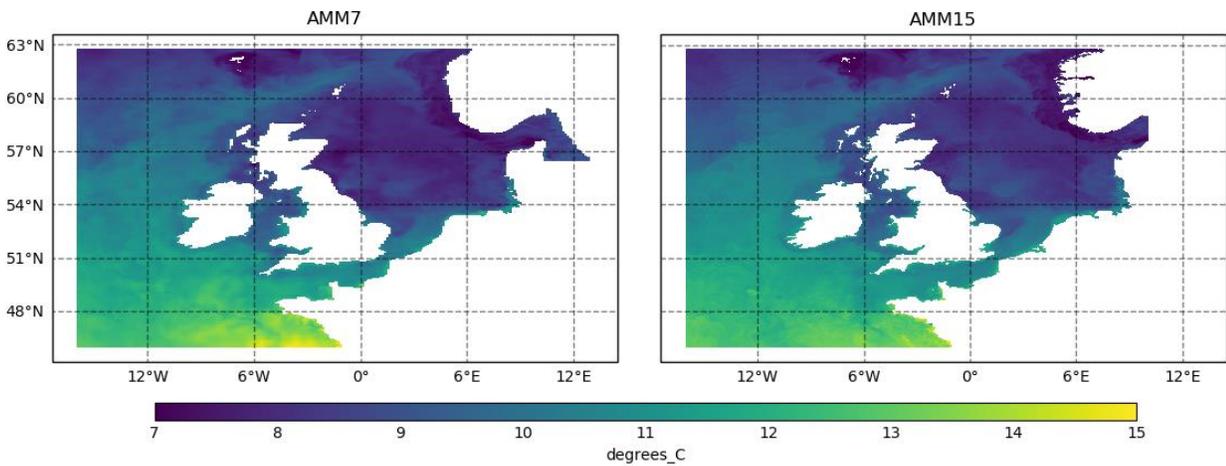
622
623 For the purposes of this assessment the 5-day daily mean sea surface potential temperature (SST)
624 forecasts (with lead times of 12, 36, 60, 84, 108 hours) were utilised for the period from January
625 to September 2019. Forecasts were compared for the co-located areas of AMM7 and AMM15.
626 Figure 1 shows the AMM7 and AMM15 co-located domain along with the land-sea mask for each
627 of the models. AMM15 has a more detailed coastline [and SST field](#) than AMM7 due to its higher
628 resolution. [When comparing two models with different resolutions it is important to know](#)
629 [whether increased detail actually translates into better forecast skill. Additionally, the](#)~~These~~
630 differences in coastline representation can have an impact on any HiRA results obtained, as will
631 be discussed in a later section.

Sea Water Potential Temperature ° C



632

Sea Water Potential Temperature



633

634 *Figure 1 - AMM7 and AMM15 co-located areas. Note the difference in the land-sea boundaries due to the different resolutions,*
635 *notably around the Scandinavian coast. [Contours show the model bathymetry at 200, 2000 and 4000 m.](#)*

636

637 It should be noted that this study is an assessment of the application of spatial methods to ocean
638 forecast data, and as such, is not meant as a full and formal assessment and evaluation of the
639 forecast skill of the AMM7 and AMM15 ocean configurations. To this end, a number of
640 considerations have had to be taken into account in order to reduce the complexity of this initial
641 study. Specifically, it was decided at an early stage to use daily mean SST temperatures, as

642 opposed to hourly instantaneous SST, as this avoided any influence of the diurnal cycle and tides
643 on any conclusions made. AMM15 and AMM7 daily means are calculated as means over 25 hours
644 to remove both the diurnal cycle and the tides. [The tidal signal is removed because the period of](#)
645 [the major tidal constituent, the semidiurnal lunar component M2, is 12 hr and 25 min \(Howarth](#)
646 [and Pugh, 1983\).](#) Daily means are also one of the variables that are available from the majority
647 of the products within the CMEMS catalogue, including reanalysis, so the application of the
648 spatial methods could be relevant in other use cases beyond those considered here. In addition,
649 there are differences in both the source and frequency of the air-sea interface forcing used in
650 both the AMM7 and AMM15 configurations which could influence the results. Most notably, the
651 AMM7 uses hourly surface pressure and 10 m winds from the Met Office Unified Model (UM),
652 whereas the AMM15 uses 3-hourly data from ECMWF.

653 2.2 Observations

654 SST observations used in the verification were downloaded from the CMEMS catalogue from the
655 product

656

- 657 • INSITU_NWS_NRT_OBSERVATIONS_013_036

658

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

666

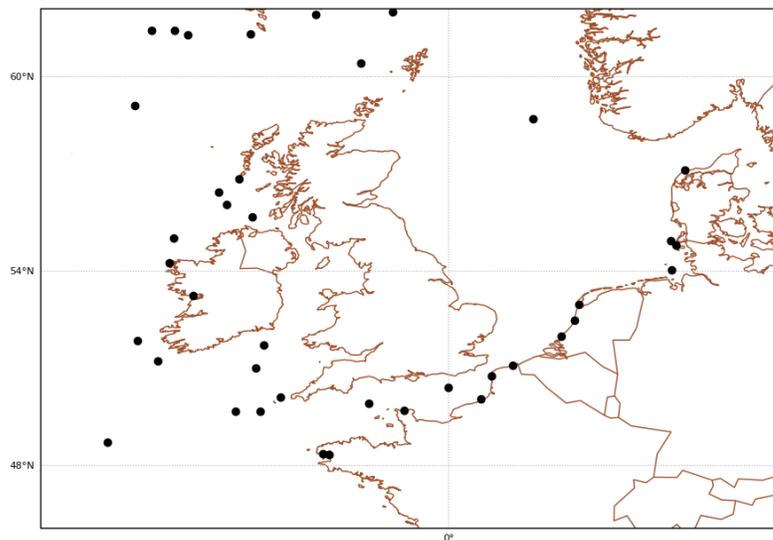
667 This study was set up to detect issues that should be considered by users when applying HiRA
668 within a routine ocean verification set-up, using a broad assessment containing as much data as

669 was available in order to understand the impact of using HiRA for ocean forecasts. Several
670 assumptions were made in this study.

671

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

682



683

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

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

688

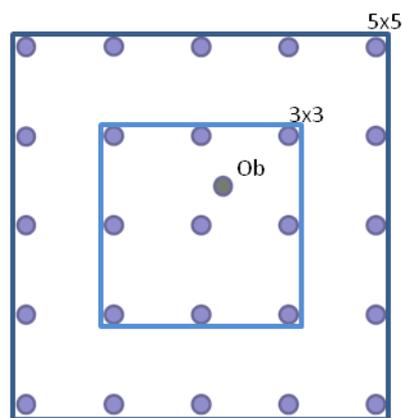
689 3. High Resolution Assessment (HiRA)

690 The HiRA framework (Mittermaier, 2014) was designed to overcome the difficulties encountered
691 in assessing the skill of high-resolution models when evaluating against point observations.
692 Traditional verification metrics such as RMSE and mean error rely on a precise matching in space
693 and time, by (typically) extracting the nearest model grid point to an observing location. The
694 method is an example of a single-observation-neighbourhood-forecast (SO-NF) approach, with
695 no smoothing. All the forecast grid points within a neighbourhood centred on an observing
696 location are treated as a pseudo ensemble, which is evaluated using well known ensemble and
697 probabilistic forecast metrics. Scores are computed for a range of (increasing) neighbourhood
698 sizes to understand the scale-error relationship. This approach assumes that the observation is
699 representative of not only its precise location but also has characteristics of the surrounding area
700 as well. WMO manual No 8 (2017) suggests that, in the atmosphere, observations can be
701 considered to be representative of an area within a 100 km radius of a land station, but this is
702 often very optimistic. The manual states further: “For small-scale or local applications the
703 considered area may have dimensions of 10 km or less.” A similar principle applies to the ocean,
704 i.e. observations can represent an area around the nominal observation location, though the
705 representative scales are likely to be very different from in the atmosphere. The representative
706 scale for an observation will also depend on local characteristics of the area, for example whether
707 the observation ~~Therefore, there~~ is on the shelf, or in open ocean or likely to be impacted by river
708 discharge.
709 There will be a limit to the useful forecast neighbourhood size which can be used when comparing
710 to a point observation. This maximum neighbourhood size will depend on the representative
711 scale, based on the representativeness of the variable under consideration. Put differently, once
712 the neighbourhoods become too big there will be forecast values in the pseudo ensemble which

713 will not be representative of the observation (and the local climatology) and any skill calculated
714 will be essentially random. [Combining results for multiple observations with very different](#)
715 [representative scales \(for example a mixture of deep ocean and coastal observations\) could](#)
716 [contaminate results, due to the forecast neighbourhood only being representative of a subset of](#)
717 [the observations. The effect of this is explored later in this paper](#)~~The scale at which~~
718 ~~representativeness is lost will vary depending on the characteristics of the variable being~~
719 ~~assessed.~~

720
721 HiRA can be based on a range of statistics, data thresholds and neighbourhood sizes in order to
722 assess a forecast model. When comparing deterministic models of different resolutions, the
723 approach is to equalise on the physical area of the neighbourhoods (i.e. having the same
724 “footprint”). By choosing sequences of neighbourhoods that provide (at least) approximate
725 equivalent neighbourhoods (in terms of area), two or more models can be fairly compared.

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



730
731 *Figure 3 - Example of forecast grid point selections for different HiRA neighbourhoods for a single observation point. A 3x3 domain*
732 *returns 9 points that represent the nearest forecast grid points in a square around the observation. A 5x5 domain encompasses*
733 *more points.*

734

735 Once neighbourhoods have been constructed, the data can be assessed using a range of well-
736 known ensemble or probabilistic scores. The choice of statistic usually depends on the
737 characteristics of the parameter being assessed. Parameters with significant thresholds can be
738 assessed using the Brier score (Brier, 1950) or the Ranked Probability Score (RPS) (Epstein, 1969),
739 i.e. assessing the ability of the forecast to correctly locate a forecast in the correct threshold
740 band. For continuous variables such as SST, the data has been assessed using the continuous
741 ranked probability score (CRPS) (Brown, 1974, Hersbach, 2000).

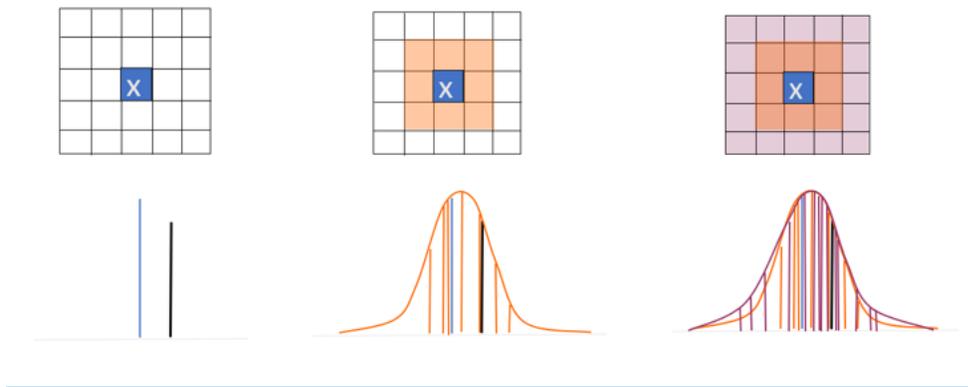
742 The CRPS is a continuous extension of the RPS. Whereas the RPS is effectively an average of a
743 user-defined set of Brier scores over a finite number of thresholds, the CRPS extends this by
744 considering an integral over all possible thresholds. It lends itself well to ensemble forecasts of
745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

749
$$CRPS = \int_{-\infty}^{\infty} [P_{fcst}(x) - P_{obs}(x)]^2 dx \quad (1)$$

750

751 Equation (1) defines the CRPS, where for a parameter x , $P_{fcst}(x)$ is the cumulative distribution of
752 the neighbourhood forecast and $P_{obs}(x)$ is the cumulative distribution of the observed value,
753 represented by a Heaviside function (see Hersbach, 2000). The CRPS is an error-based score
754 where a perfect forecast has a value of zero. It measures the difference between two cumulative
755 distributions, a forecast distribution formed by ranking the (in this case quasi) -ensemble
756 members represented by the forecast values in the neighbourhood, and a step function
757 describing the observed state. To use an ensemble, HiRA makes the assumption that all grid
758 points within a neighbourhood are equi-probable outcomes at the observing location. Therefore,
759 aside from the observation representativeness limit, as the neighbourhood sizes increase, this
760 assumption of equi-probability will break down as well, and scores become random. Care must

761 therefore be taken to decide whether a particular neighbourhood size is appropriately
762 representative. This decision will be based on the length scales appropriate for a variable as well
763 as the resolution of the forecast model being assessed. [Figure 4 shows a schematic of how](#)
764 [different neighbourhood sizes contribute towards constructing forecast probability density](#)
765 [functions around a single observation.](#)



766
767 [Figure 4 – Example of how different forecast neighbourhood sizes would contribute to generation of a probability density function](#)
768 [around an observation \(denoted by x\). The larger the neighbourhood, the better described the pdf, though potentially at the](#)
769 [expense of larger spread. Where a forecast point is invalid within the forecast neighbourhood then that site is rejected from the](#)
770 [calculations for that neighbourhood size.](#)

771
772 AMM7 and AMM15 resolve different length scale of motion, due to their horizontal resolution.
773 This should be taken into account when assessing the results of different neighbourhood sizes.
774 Both models can resolve the large barotropic scale (~200 km) and the shorter baroclinic scale off
775 the shelf, in deep water. On the continental shelf, only the resolution of ~1.5 km of AMM15,
776 permits motions at the smallest baroclinic scale since the first baroclinic Rossby radius is of order
777 of 4 km (O’Dea et al., 2012). AMM15 represents a step change in representing the eddy dynamics
778 variability on the continental shelf. This difference has an impact also on the data assimilation
779 scheme, where two horizontal correlation length scales (Mirouze et al., 2016) are used to
780 represent large and small scales of ocean variability. The long length scale is 100 km while the
781 short correlation length scale aims to account for internal ocean processes variability,
782 characterized by the Rossby radius of deformation. Computational requirements restrict the

783 short length scale to be at least 3 model grid points, 4.5 km and 21 km respectively for AMM15
784 and AMM7 (Tonani et al., 2019). Although AMM15 resolves smaller scale processes, comparing
785 AMM7 and AMM15 in neighbourhood sizes between the AMM7 resolution and multiples of this
786 resolution will address processes that should be accounted for in both models.

787

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

792

793 4. Model Evaluation Tools (MET)

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

798

799 5. Equivalent neighbourhoods and equalisation

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

806

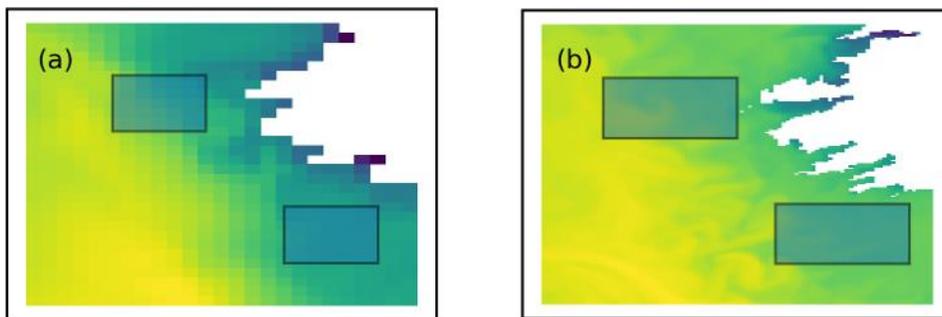
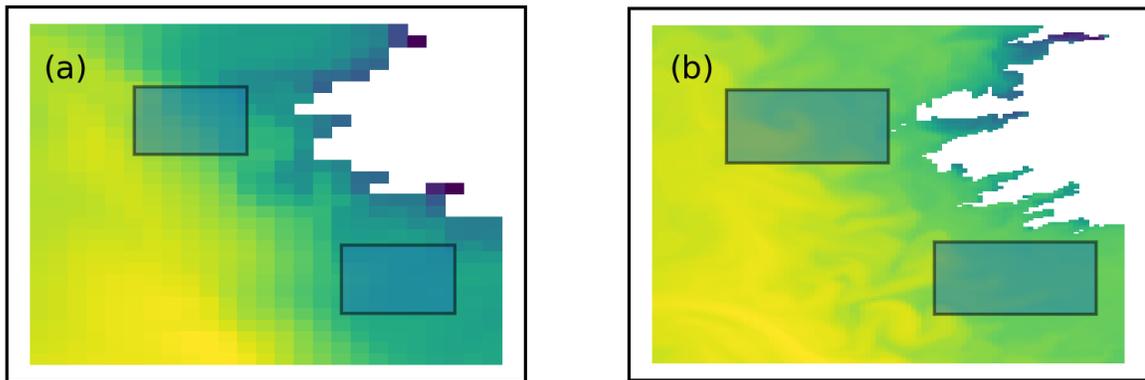
807 The two model configurations used in this assessment are provided on standard latitude-
808 longitude grids via the CMEMS catalogue. The AMM7 and AMM15 configurations are stated to

809 have resolutions approximating 7 km and 1.5 km respectively. Thus, equivalent neighbourhoods
810 should simply be a case of matching neighbourhoods with similar spatial distances. In fact, the
811 AMM15 is originally run on a rotated latitude-longitude grid where the resolution is closely
812 approximated by 1.5 km and subsequently provided to the CMEMS catalogue on the standard
813 latitude-longitude grid. Once the grid has been transformed to a regular latitude-longitude grid
814 the 1.5 km nominal spatial resolution is not as accurate. This is particularly important when
815 neighbourhood sizes become larger, since any error in the approximation of the resolution will
816 become multiplied as the number of points being used increases.

817

818 Additionally, the two model configurations do not have the same aspect ratio of grid points.
819 AMM7 has a longitudinal resolution of $\sim 0.11^\circ$ and a latitudinal resolution of $\sim 0.066^\circ$ (a ratio of
820 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).
821 HiRA neighbourhoods typically contain the same number of grid-points [in the zonal](#)~~vertically~~ and
822 [meridional directions](#)~~horizontally~~ which will lead to discrepancies in the area selected when
823 comparing models with different grid aspect ratios, depending on whether the comparison is
824 based on neighbourhoods with a similar longitudinal or similar latitudinal size. This difference will
825 scale as the neighbourhood size increases as shown in Fig. 4 [and Table 2](#). The onus is therefore
826 on the user to understand any difference in grid structure, and therefore [within the](#) HiRA
827 neighbourhoods, between models being compared and to allow for this when comparing
828 equivalent neighbourhoods.

829



(c)	AMM7		AMM15		Size (E-W)	
	Name	Total points	Shape	Total points	Shape	Degrees
NB1	1	1x1	25	5x5	0.11	7
NB2	9	3x3	121	11x11	0.33	21
NB3	25	5x5	361	19x19	0.55	35
NB4	49	7x7	625	25x25	0.77	49
NB5	81	9x9	1089	33x33	0.99	63

833 *Figure 5 - Similar neighbourhood sizes for a 49 km neighbourhood using the approximate resolutions (7 km and 1.5 km) with a)*
 834 *AMM7 with a 7x7 neighbourhood (NB4), b) AMM15 with a 33x33 neighbourhood (NB5) and c) details of equivalent*
 835 *neighbourhood sizes and naming conventions, with scales relating to AMM7. Whilst the neighbourhoods are similar sizes in the*
 836 *latitudinal direction, the AMM15 neighbourhood is sampling a much significantly larger area due to different scales in the*
 837 *longitudinal direction. This means that a comparison with a 25x25 AMM15 neighbourhood is more appropriate.*

838 [Table 2 - Details of equivalent neighbourhoods used when comparing AMM7 and AMM15.](#)

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

839

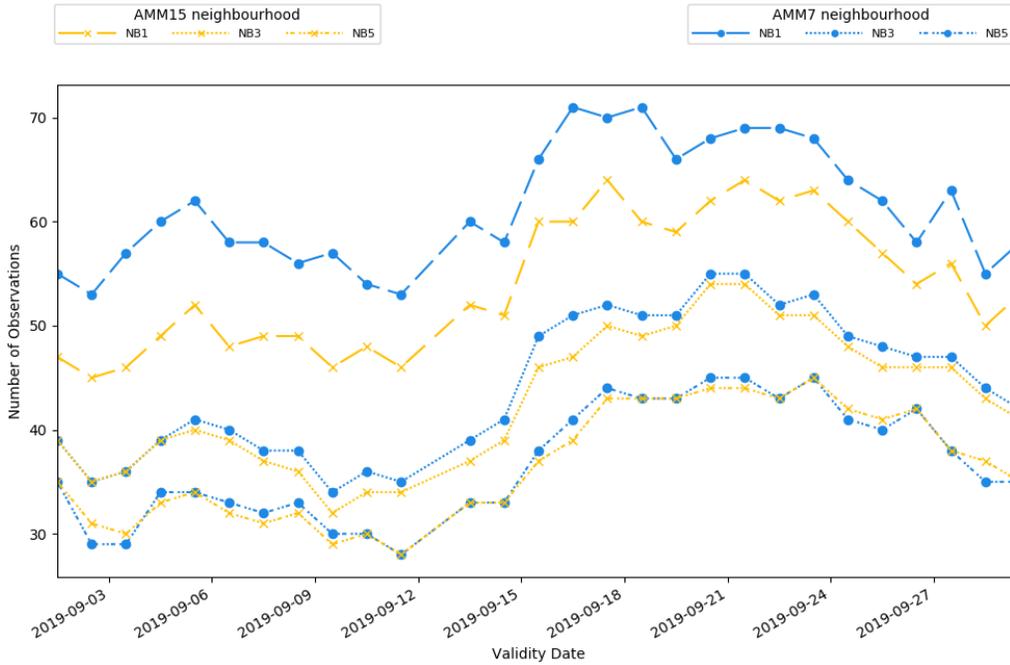
840 For this study we have matched neighbourhoods between model configurations based on their
 841 longitudinal size. The equivalent neighbourhoods used to show similar areas within the two
 842 configurations are indicated in [Table 2](#) ~~Fig. 4e~~ along with the bar style and naming convention
 843 used throughout.

844

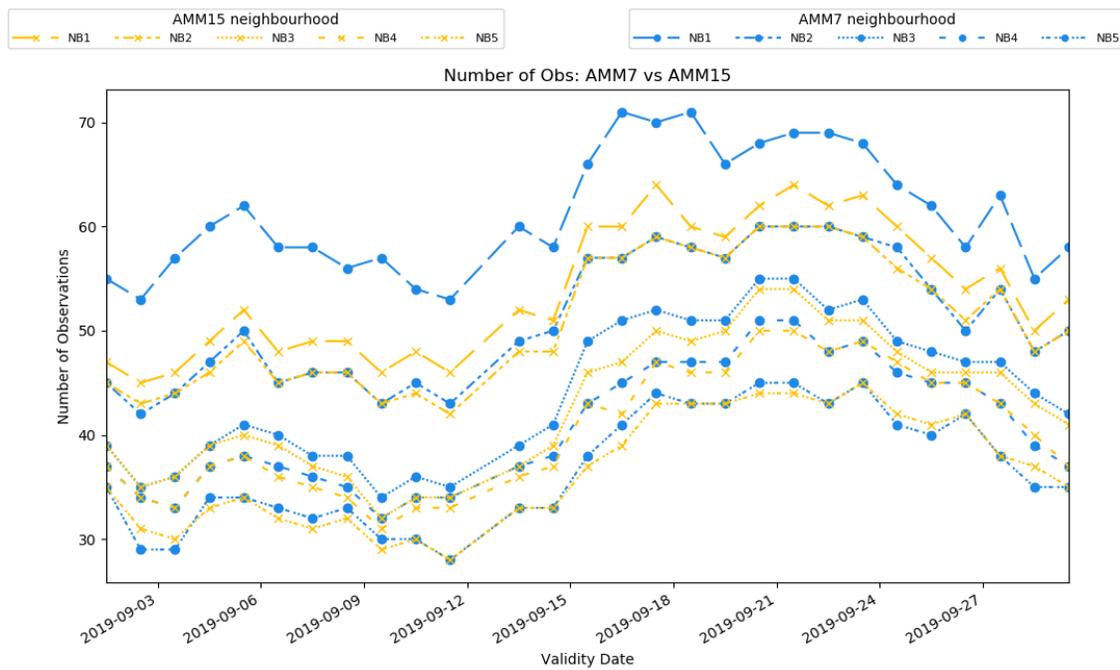
845 For ocean applications there are other aspects of the processing to be aware of when using
 846 neighbourhood methods. This is mainly related to the presence of coastlines and how their
 847 representation changes resolution (as defined by the land-sea mask) and the treatment of
 848 observations within HiRA neighbourhoods. Figure [54](#) illustrates the contrasting land-sea
 849 boundaries due to the different resolutions of the two configurations. When calculating HiRA
 850 neighbourhood values, all forecast values in the specific neighbourhood around an observation
 851 must be present for a score to be calculated. [If any forecast points within a neighbourhood
 852 contain missing data then that observation at that neighbourhood size is rejected.](#) This is to
 853 ensure that the resolution of the “ensemble”, which is defined or determined by the number of
 854 members, remains the same. For typical atmospheric fields such as screen temperature this is
 855 not an issue, but with parameters that have physical boundaries (coastlines), such as SST, there
 856 will be discontinuities in the forecast field that depend on the location of the land-sea boundary.
 857 For coastal observations, this means that as the neighbourhood size increases, it is more likely
 858 [that an observation will](#) ~~to~~ be rejected from the comparison due to missing data. Even at the grid
 859 scale, the nearest model grid point to an observation may not be a sea point. In addition, different

860 land-sea borders between models mean that potentially some observations will be rejected from
861 one model comparison but will be retained in the other [because of missing forecast points within](#)
862 [their respective neighbourhoods.](#) Care should be taken when implementing HIRA to check the
863 observations available to each model configuration when assessing the results and make a
864 judgement as to whether the differences are important.

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



871



872

873 Figure 6- Number of observation sites within NB1, NB3 and NB5for each neighbourhood size for AMM15 and AMM7. Numbers
 874 are those used during September 2019 but represent typical total observations during a month. Matching line styles represent
 875 equivalent neighbourhoods.

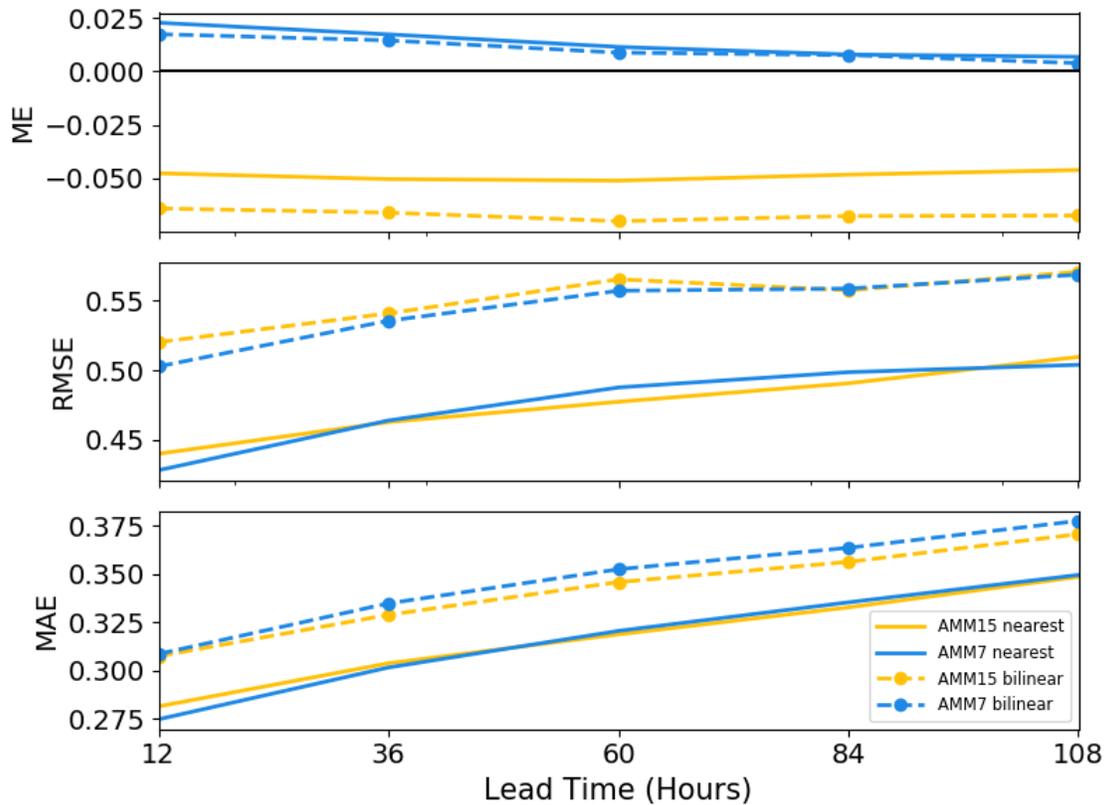
876

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

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

897

898 6. Results



899

900 *Figure 7 - Verification results using a typical statistics approach for January – September 2019. Mean error (top), root mean square*
 901 *error (middle) and mean absolute error (bottom) results are shown for the two model configurations. Two methods of matching*
 902 *forecast to observations points have been used; a nearest neighbor approach (solid) representing the single grid point results from*
 903 *HiRA, and a bilinear interpolation approach (dashed) more typically used in operational ocean verification.*

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

911 It is noted that the two methods of matching forecasts to observation locations give quite
 912 different results. For the mean error, the impact of moving from a single grid point approach to
 913 a bilinear interpolation method appears to be minor for the AMM7 model, but is more severe for

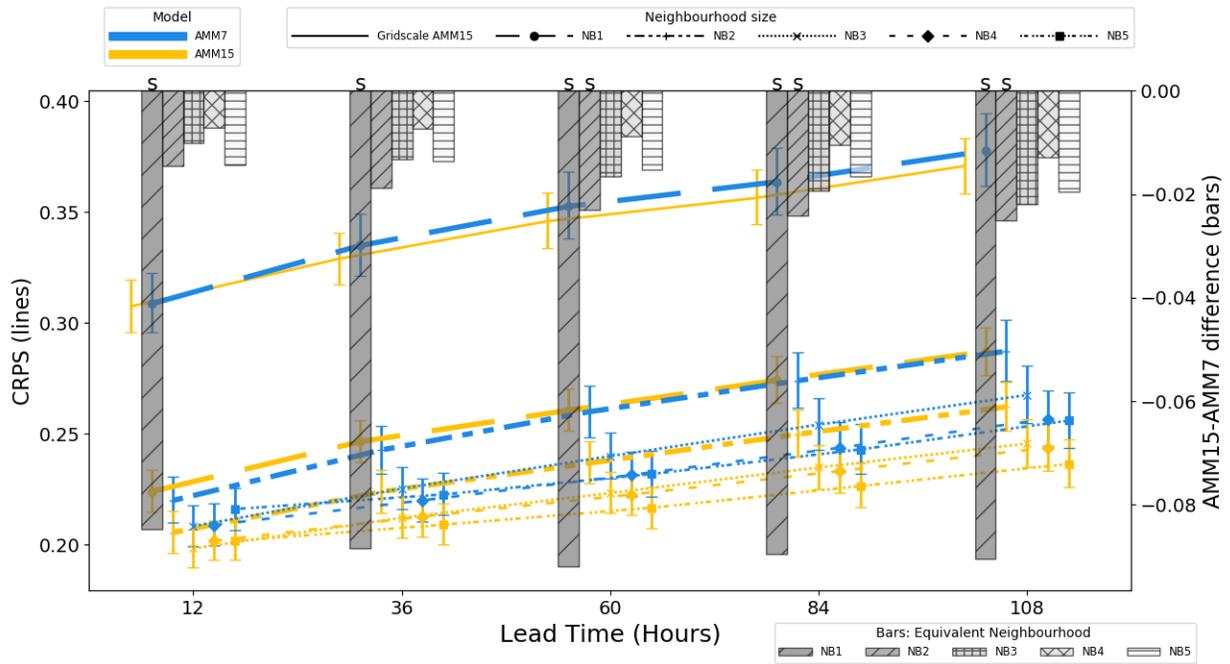
914 the AMM15, resulting in a larger error across all lead times. For the RMSE the picture is more
915 mixed, generally suggesting that the AMM7 forecasts are better when using a bilinear
916 interpolation method but giving no clear overall steer when the nearest grid point is used.
917 However, the impact of taking a bilinear approach results in much higher gross errors across all
918 lead times when compared to the nearest grid point approach.

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

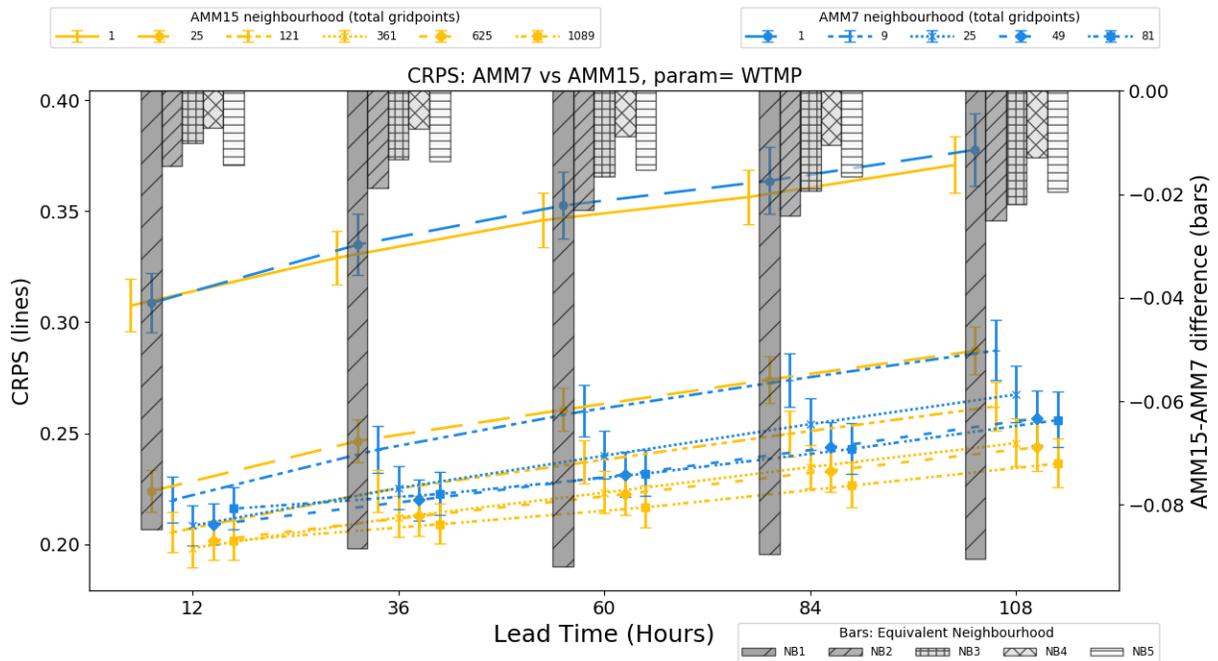
925 Based on the interpolated RMSE results in Fig. 6 it would be hard to conclude that there was a
926 significant benefit to using high-resolution ocean models for forecasting SSTs. This is where the
927 HiRA framework can be applied. It can be used to provide more information, which can better
928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

929

930



931



932

933

934

935

936

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

937 Figure [87](#) shows the results for AMM7 and AMM15 for the period January - September 2019
938 using the HiRA framework with the CRPS. The lines on the plot show the CRPS for the two model
939 configurations for different neighbourhood sizes, each plotted against lead-time. Similar line
940 styles are used to represent equivalent neighbourhood sizes. Confidence intervals have been
941 generated by applying a bootstrap with replacement method, using 10000 samples, to the
942 domain-averaged CRPS (e.g. Efron and Tibshirani, 1993). The error bars represent the 95 %
943 confidence level. The results for the single grid-point show the MAE and are the same as would
944 be obtained using a traditional (precise) matching. In the case of CRPS, where a lower score is
945 better, we see that AMM15 is better than AMM7, though not significantly so, except at shorter
946 lead-times where there is little difference.

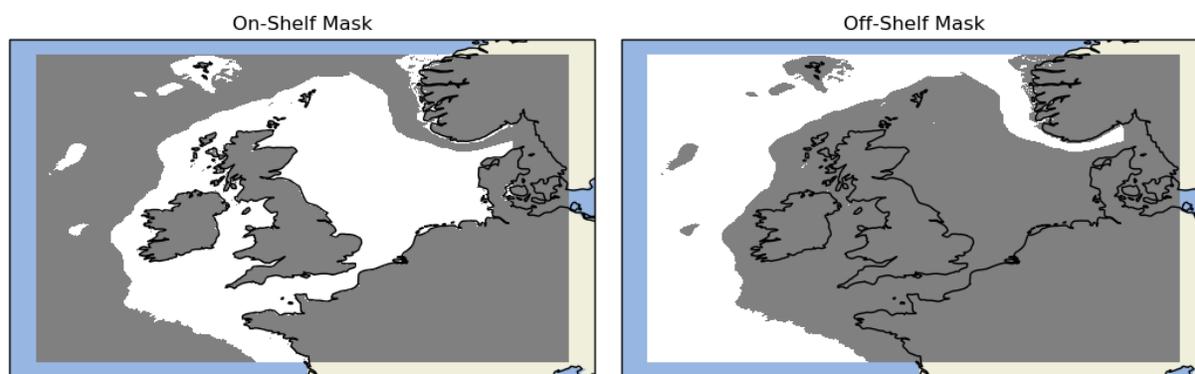
947 The differences at equivalent neighbourhood sizes are displayed as a bar plot on the same figure,
948 with scores referenced with respect to the right-hand axis. Line markers and error bars have been
949 offset to aid visualization, such that results for equivalent neighbourhoods are displayed in the
950 same vertical column as the difference indicated by the barplot. The details of the equivalent
951 neighbourhood sizes are presented in [Table 2, Fig. 4c](#). Since a lower CRPS score is better, a
952 positively orientated (upwards) bar implies AMM7 is better, whilst a negatively orientated
953 (downwards) bar means AMM15 is better.

954 As [indicated defined](#) in [Table 2, Fig. 4c](#) NB1 compares the single grid-point results of AMM7 with
955 a 25-member pseudo-ensemble constructed from a 5x5 AMM15 neighbourhood. Given the
956 different resolutions of the two configurations, these two neighbourhoods represent similar
957 physical areas from each model domain, with AMM7 only represented by a single forecast value
958 for each observation, but AMM15 represented by 25 values cover the same area, and as such
959 potentially better able to represent small-scale variability within that area.

960 At this equivalent scale the AMM15 results are markedly better than AMM7, with lower errors,
961 suggesting that overall the AMM15 neighbourhood better represents the variation around the
962 observation than the coarser single grid point of AMM7. At the next set of equivalent
963 neighbourhoods (NB2), the gap between the two configurations has closed, but AMM15 is still
964 consistently better than AMM7 as lead time increases. Above this scale the neighbourhood

965 values tend towards similarity, and then start to diverge again suggesting that the representative
966 scale of the neighbourhoods has been reached and that errors are essentially random.

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



975
976 *Figure 9 - On-shelf and off-shelf masking regions within the co-located AMM7 and AMM15 domain (data within the grey areas is*
977 *masked).*

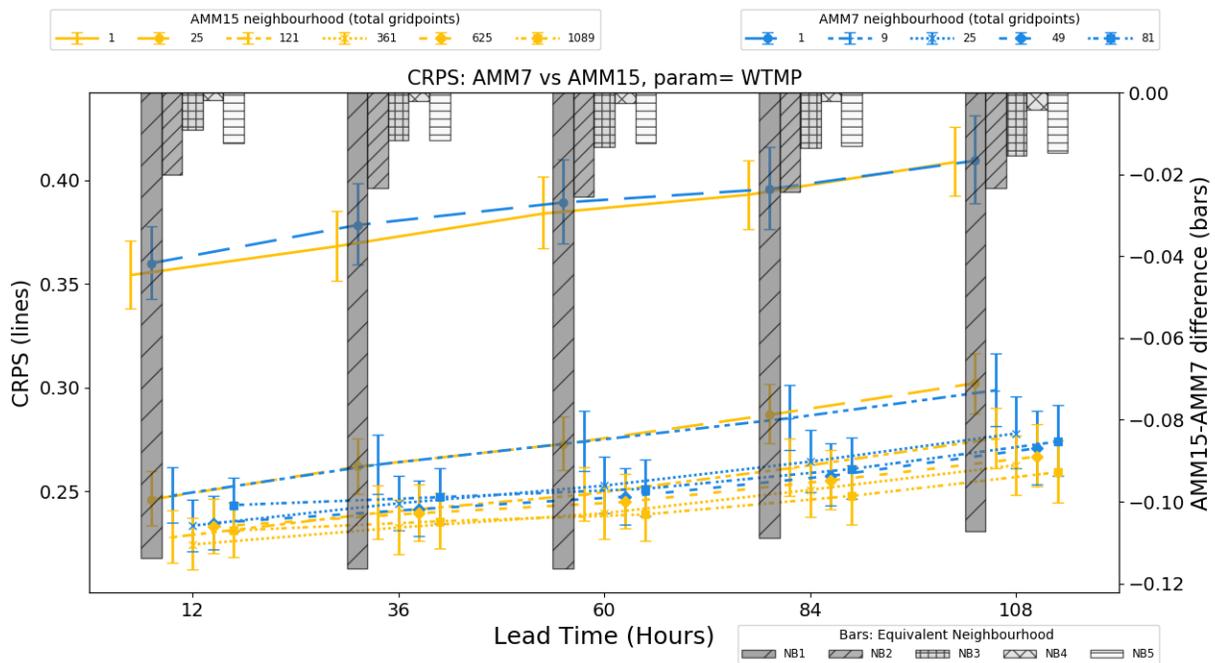
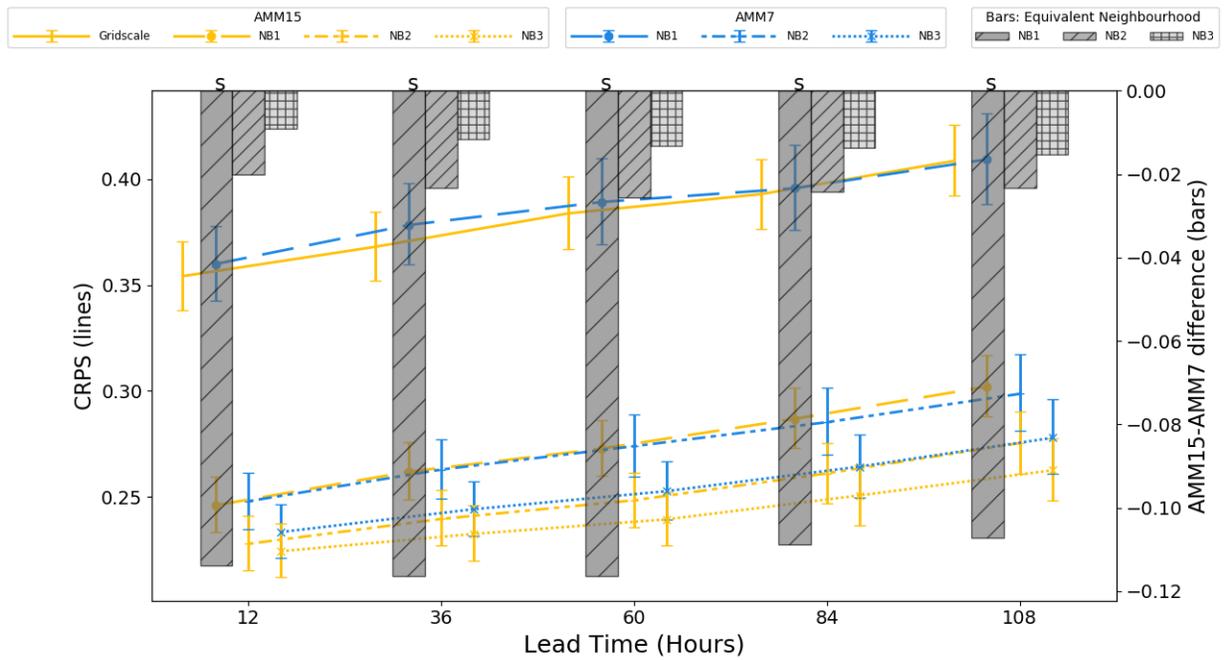


Figure 10- Summary of on-shelf CRPS (left axis, lines) and CRPS difference (right axis, bars) for the period January 2019 to September 2019 for AMM7 and AMM15 models at different neighbourhood sizes. Error bars represent 95 % confidence values obtained from 10000 samples using bootstrap with replacement.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

985 On-shelf results (Fig. [109](#)) show that at the grid scale the results for both AMM7 and AMM15 are
986 worse for this sub-domain. This could be explained by both the complexity of processes (tides,
987 friction, river mixing, topographical effects, etc.), and the small dynamical scales associated with
988 shallow waters on the shelf (Holt et al., 2017).

989

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

996

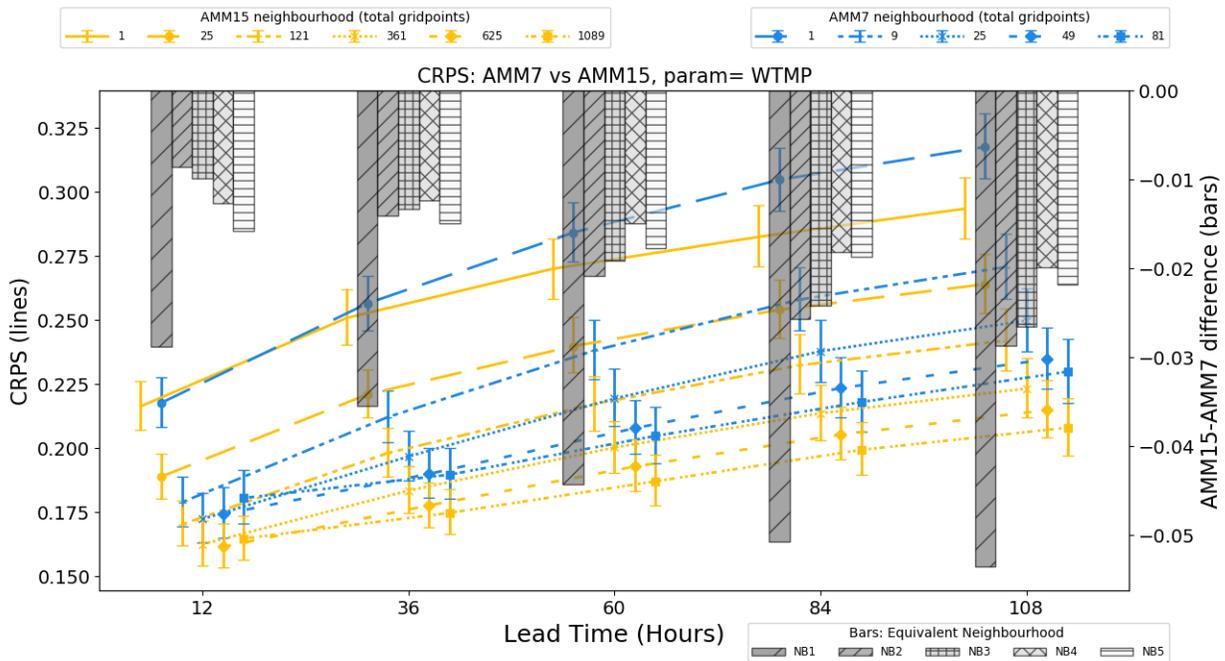
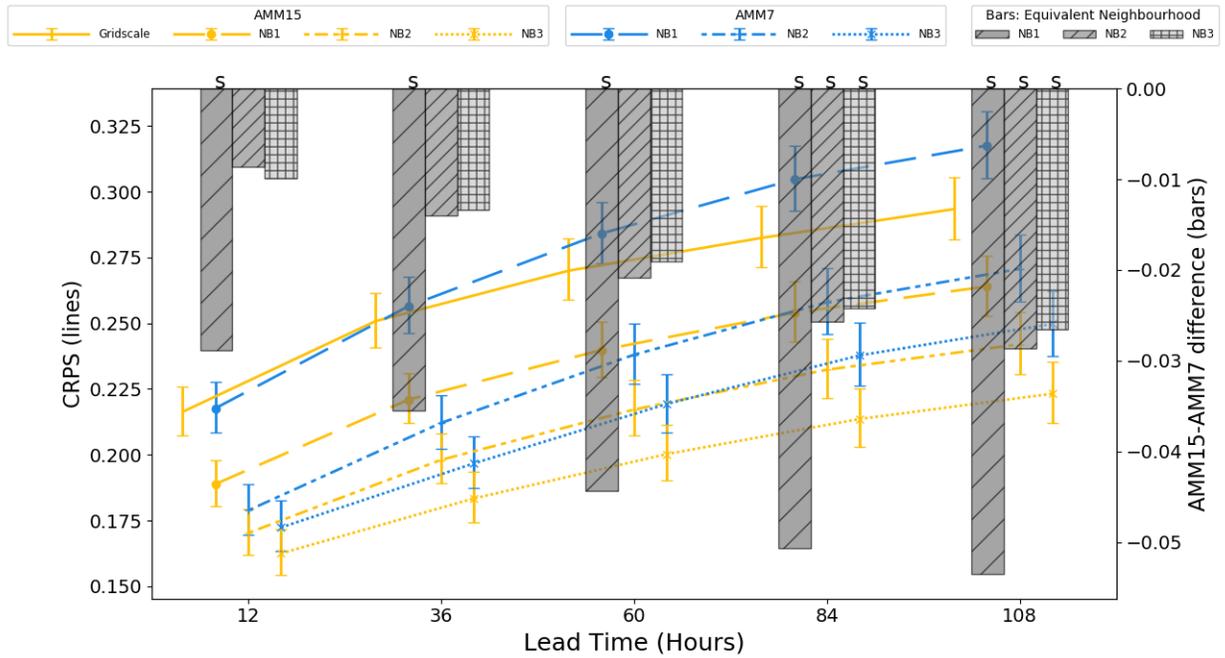


Figure 11 – Summary of off-shelf CRPS (left axis, lines) and CRPS difference (right axis, bars) for the period January 2019 to September 2019 for AMM7 and AMM15 models at different neighbourhood sizes. Error bars represent 95 % confidence values obtained from 10000 samples using bootstrap with replacement.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

1005 For off-shelf results (Fig. [1110](#)), the CRPS is much better (smaller error), at both the grid scale and
1006 for HiRA neighbourhoods, suggesting that both configurations are better at forecasting these
1007 deep ocean SSTs (or that it is easier to do so). There is still an improvement in CRPS when going
1008 from the grid scale (single grid box) to neighbourhoods, but the value of that change is much
1009 smaller than for the on-shelf sub-domain. When comparing equivalent neighbourhoods, the
1010 AMM15 still gives consistently better results (smaller errors) and appears to improve over AMM7
1011 as lead time increases in contrast to the on-shelf results.

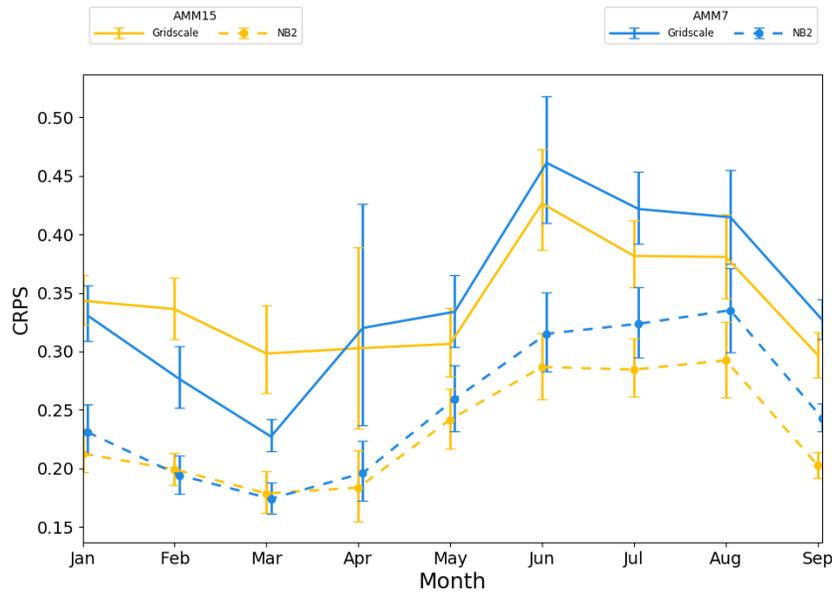
1012 It is likely that the neighbourhood at which we lose representativity will be larger for the deeper
1013 ocean than the shelf area because of the larger scale of dynamical processes in deep water. When
1014 choosing an optimum neighbourhood to use for assessment, care should be taken to check
1015 whether there are different representativity levels in the data (such as here for on-shelf and off-
1016 shelf) and pragmatically choose the smaller of those equivalent neighbourhoods when looking at
1017 data combining the different representativity levels.

1018 Overall, for the period January-September 2019, the AMM15 demonstrates a lower (better) CRPS
1019 than AMM7 when looking at the HiRA neighbourhoods. However, this also appears to be true at
1020 the grid scale over the assessment period. One of the aspects that HiRA is trying to provide
1021 additional information about is whether higher resolution models can demonstrate improvement
1022 over coarser models against a perception that the coarser models score better in standard
1023 verification forecast assessments. Assessed over the whole period, this initial premise does not
1024 appear to hold true, therefore a [deeper](#)~~closer~~ look at the data is required [to assess whether this
1025 signal is consistent within shorter time periods, or whether there are underlying periods
1026 contributing significant and contrasting results to the whole-period aggregate.](#) -

1027 Figure 12 shows a monthly breakdown of the grid scale and the NB2 HiRA neighbourhood scores
1028 at T+60. This shows the underlying monthly variability not immediately apparent in the whole-
1029 period plots. Notably for the January to March period, AMM7 outperforms AMM15 at the grid
1030 scale. With the introduction of HiRA neighbourhoods, AMM7 still performs better for February
1031 and March but the difference between the models is significantly reduced. For these monthly
1032 timeseries the error bars increase in size relative to the summary plots (e.g. Fig [87](#)) due to the

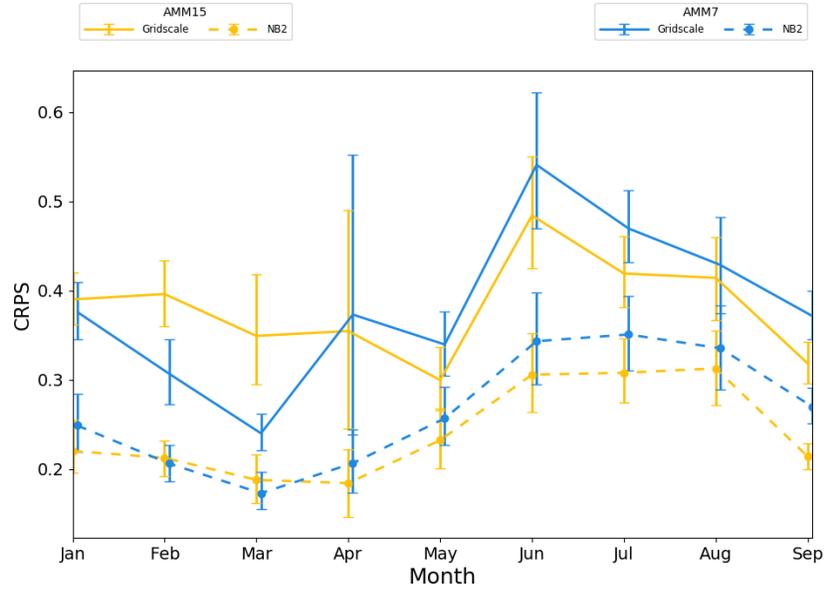
1033 reduction in data available. The sample size will have an impact on the error bars as the smaller
 1034 the sample, the less representative of the true population the data is likely to be. April in
 1035 particular ~~contained~~contains several days of missing forecast data, leading to a reduction in
 1036 sample size and corresponding increase in error bar size, whilst during May there was a period
 1037 with reduced numbers of observations.

1038



1039

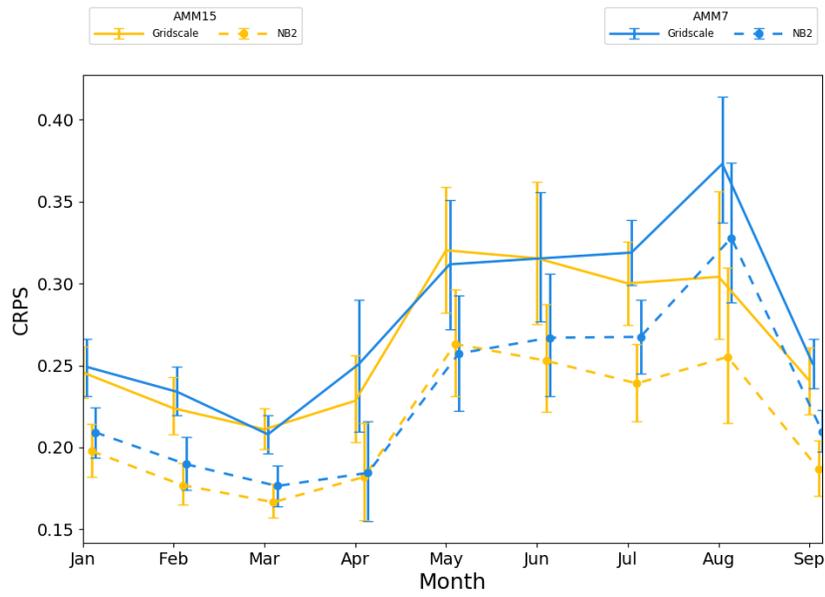
1040 *Figure 12 – Monthly time series of whole-domain CRPS scores for grid scale (solid line) and NB2 neighbourhood (dashes) for T+60*
 1041 *forecasts. Error bars represent 95 % confidence values obtained from 10000 samples using bootstrap with replacement. Error bars*
 1042 *have been staggered in the x-direction to aid clarity.*



1043

1044 *Figure 13 - On-shelf monthly time series of CRPS. Error bars represent 95 % confidence values obtained from 10000 samples using*
 1045 *bootstrap with replacement.*

1046



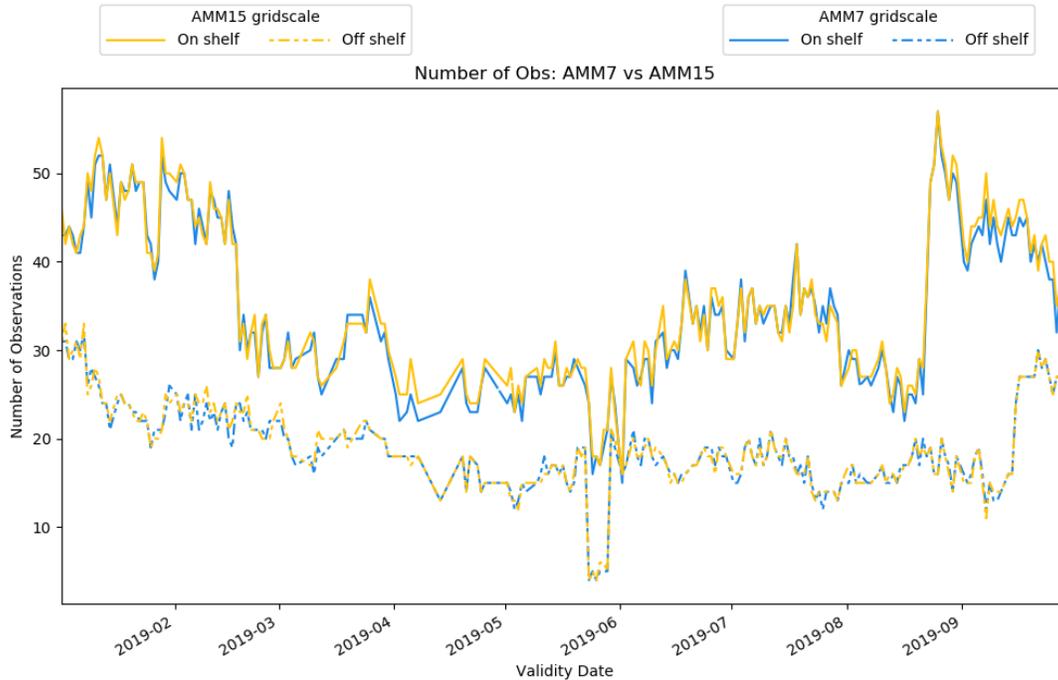
1047

1048 *Figure 14 - Off-shelf monthly time series of CRPS. Error bars represent 95 % confidence values obtained from 10000 samples using*
 1049 *bootstrap with replacement.*

1050

1051 The same pattern is present for the on-shelf sub-domain (Fig. [1312](#)), where what appears to be
1052 a significant benefit for the AMM7 during February and March is less clear-cut at the NB2
1053 neighbourhood. For the off-shelf sub-domain (Fig. [1413](#)), differences between the two
1054 configurations at the grid scale are mainly apparent during the summer months. At the NB2 scale,
1055 the AMM15 [potentially](#) demonstrates more benefit than AMM7 except for April and May, where
1056 the two show similar results. [There is a balance to be struck in this conclusion as the differences](#)
1057 [between the two models are rarely greater than the 95 % error bars. This in itself does not mean](#)
1058 [that the results are not significant. However, care should be taken when interpreting such a result](#)
1059 [as a statistical conclusion rather than broad guidance as to model performance. Attempts to](#)
1060 [reduce the error bar size, such as increasing the number of observations, or number of times](#)
1061 [within the period would aid this interpretation.](#)

1062 One noticeable aspect of the time series plots is that the whole-domain plot is heavily influenced
1063 by the on-shelf results. This is due to the difference in observation numbers as shown in Fig. [1514](#),
1064 with the on-shelf domain having more observations overall, sometimes significantly more, for
1065 example during January or mid-late August. For the overall domain, the on-shelf observations
1066 will contribute more to the overall score and hence the underlying off-shelf signal will tend to be
1067 masked. This is an indication of why verification is more useful when done over smaller, more
1068 homogeneous sub-regions, rather than verifying everything together, with the caveat that
1069 sample sizes are large enough, since underlying signals can be swamped by dominant error types.



1070

1071 *Figure 15 - Number of grid scale observations for the on and off-shelf domains.*

1072

1073 **7. Discussion and Conclusions**

1074 In this study, the HiRA framework has been applied to SST forecasts from two ocean models with
 1075 different resolutions. This enables a different view of the forecast errors than obtained using
 1076 traditional (precise) grid scale matching against ocean observations. Particularly it enables us to
 1077 demonstrate the additional value of high-resolution model. When considered more
 1078 appropriately high-resolution models (with the ability to forecast small-scale detail) have lower
 1079 errors when compared to the smoother forecasts provided by a coarser-resolution model.

1080 The HiRA framework was intended to address the question ‘Does moving to higher resolution
 1081 add value?’ This study has identified and highlighted aspects that need to be considered when
 1082 setting up such an assessment. Prior to this study, routine verification statistics typically showed
 1083 that coarser resolution models had equivalent or more skill than higher resolution models (e.g.
 1084 Mass et al., 2002, Tonani et al., 2019). During the period January to September 2019, grid scale

1085 verification within this assessment showed that the coarser-resolution AMM7 often
1086 demonstrated lower errors than the AMM15.

1087 HiRA neighbourhoods were applied and the data then assessed using the CRPS, showing a large
1088 reduction (improvement) in errors for AMM15 when going from a grid scale, point-based
1089 verification assessment to a neighbourhood, ensemble approach. When applying an equivalent-
1090 sized neighbourhood to both configurations, AMM15 typically demonstrated lower (better)
1091 scores. These scores were in turn broken down into off-shelf and on-shelf sub-domains and
1092 showed that the different physical processes in these areas affected the results. [Forecast
1093 verification studies tailored for the coastal/shelf areas are needed to properly understand the
1094 forecast skills in areas with high complexity and fast evolving dynamics.](#)

1095 When constructing HiRA neighbourhoods the spatial scales that are appropriate for the
1096 parameter must be considered carefully. This often means running at several neighbourhood
1097 sizes and determining where the scores no longer seem physically representative. When
1098 comparing models, care should be taken to construct neighbourhood sizes that are similarly sized
1099 spatially, the details of the neighbourhood sizes will depend on the structure and resolution of
1100 the model grid.

1101 Treatment of observations is also important in any verification set-up. For this study, the fact that
1102 there are different numbers of observations present at each neighbourhood scale (as
1103 observations are rejected due to land contamination) means that there is never an optimally
1104 equalized data set (i.e. the same observations for all models and for all neighbourhood sizes). It
1105 also means that comparison of the different neighbourhood results from a single model is ill
1106 advised, in this case, as the observations numbers can be very different, and therefore the model
1107 forecast is being sampled at different locations. Despite this, observation numbers should be
1108 similar when looking at matched spatially sized neighbourhoods from different models if results
1109 are to be compared. One of the main constraints identified through this work is both the sparsity
1110 and geographical distribution of observations throughout the North West Shelf domain, with
1111 several viable locations rejected during the HiRA processing due to their proximity to coastlines.

1112 The purest assessment, in terms of observations, would involve a fixed set of observations,
1113 equalized across both model configurations and all neighbourhoods at every time. This would
1114 remove the variation in observation numbers seen as neighbourhood sizes increase as well as
1115 those seen between the two models and give a clean comparison between two models.

1116 Care should be taken when applying strict equalization rules as this could result in only a small
1117 number of observations being used. The total number of observations used should be large
1118 enough to ensure that the sample is large enough to produce robust results and satisfy rules for
1119 statistical significance. Equalisation rules could also unfairly affect the spatial sampling of the
1120 verification domain. For example, in this study coastal observations would be affected more than
1121 deep ocean observations if neighbourhood equalization were applied, due to the proximity of
1122 the coast.

1123 To a lesser extent, the variation in observation numbers on a day-to-day timescale also has an
1124 impact on any results and could mean that incorrect importance is attributed to certain results,
1125 which are simply due to fluctuations in observation numbers.

1126 The fact that the errors can be reduced through the use of neighbourhoods shows that the ocean
1127 and the atmosphere have similarities in the way the forecasts behave as a function of resolution.
1128 This study did not consider the concept of skill, which incorporates the performance of the
1129 forecast relative to a pre-defined benchmark. For the ocean the choice of reference needs to be
1130 considered. This could be the subject of further work.

1131 To our knowledge, this work is the first attempt to use neighbourhood techniques to assess ocean
1132 models. The promising results showing reductions in errors of the finer resolution configuration
1133 warrant further work. We see a number of directions the current study could be extended.

1134 The study was conducted on daily output which should be appropriate to address eddy mesoscale
1135 variability, but observations are distributed at hourly resolution, and so the next logical step
1136 would be to assess the hourly forecasts against the hourly observation and see how this impacted
1137 the results. This will increase the sample size, if all hourly observations were considered together.
1138 However, it is impossible to speculate on whether considering hourly forecasts would lead to
1139 more noisy statistics, counteracting the larger sample size.

1140 [This assessment only looked at SST for this initial examination.](#) Consideration of other ocean
1141 variables would also be of interest, including looking at derived diagnostics such as mixed layer
1142 depth, but the sparsity of observations available for some variables may limit the case studies
1143 available. HiRA as a framework is not remaining static. Enhancements to introduce non-regular
1144 flow-dependent neighbourhoods are planned and may be of benefit to ocean applications in the
1145 future. Finally, an advantage of using the HiRA framework is that results obtained from
1146 deterministic ocean models could also be compared against results from ensemble models when
1147 these become available for ocean applications.

1148

1149 8. References

1150 [Aznar, R., Sotillo, M., Cailleau, S., Lorente, P., Levier, B., Amo-Baladrón, A., Reffray, G. and Alvarez Fanjul,](#)
1151 [E.: Strengths and weaknesses of the CMEMS forecasted and reanalyzed solutions for the Iberia-Biscay-](#)
1152 [Ireland \(IBI\) waters. J. Marine. Syst., 159, <https://doi.org/10.1016/j.jmarsys.2016.02.007>, 2016.](#)

1153 Brassington, G.: Forecast Errors, Goodness, and Verification in Ocean Forecasting, *J. Marine Res.*, 75,
1154 403-433, <https://doi.org/10.1357/002224017821836851>, 2017.

1155 Brier, G. W.: Verification of Forecasts Expressed in Terms of Probability, *Mon. Wea. Rev.*, 78, 1-3,
1156 [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.

1157 Brown, T. A.: Admissible scoring systems for continuous distributions, Santa Monica, CA: RAND
1158 Corporation, available at <https://www.rand.org/pubs/papers/P5235.html>, 1974.

1159 Casati, B., Ross, G. and Stephenson, D. B.: A new intensity-scale approach for the verification of spatial
1160 precipitation forecasts, *Met. Apps.*, 11, 141-154, <https://doi.org/10.1017/S1350482704001239>, 2004.

1161 Davis, C., Brown, B. and Bullock, R.: Object-Based Verification of Precipitation Forecasts. Part I:
1162 Methodology and Application to Mesoscale Rain Areas, *Mon. Wea. Rev.*, **134**, 1772–1784,
1163 <https://doi.org/10.1175/MWR3145.1>, 2006.

1164 Dorninger, M., Gilleland, E., Casati, B., Mittermaier, M. P., Ebert, E. E., Brown, B. G., and Wilson, L. J.: The
1165 Setup of the MesoVICT Project, *Bull. Amer. Meteor. Soc.*, 99, 1887–1906,
1166 <https://doi.org/10.1175/BAMS-D-17-0164.1>, 2008.

1167 Ebert, E. E.: Fuzzy verification of high-resolution gridded forecasts: a review and proposed framework,
1168 Met. Apps, 15, 51-64, <https://doi.org/10.1002/met.25>, 2008.

1169 Efron, B. and Tibshirani, R.: Bootstrap methods for standard errors, confidence intervals, and other
1170 measures of statistical accuracy, Statistical Science, 1, 54-77, 1986.

1171 Epstein, E. S.: A Scoring System for Probability Forecasts of Ranked Categories, J. Appl. Meteor., 8, 985–
1172 987, 1969.

1173 Gilleland, E., Ahijevych, D., Brown, B. G., Casati, B., and Ebert, E. E.: Intercomparison of Spatial Forecast
1174 Verification Methods, Wea. Forecasting, 24, 1416–1430, <https://doi.org/10.1175/2009WAF2222269.1>,
1175 2009.

1176 Goessling, H. F., Tietsche, S., Day, J. J., Hawkins, E., and Jung, T. : Predictability of the Arctic sea ice edge,
1177 Geophys. Res. Lett., 43, 1642– 1650, doi:[10.1002/2015GL067232](https://doi.org/10.1002/2015GL067232), 2016.

1178 Graham, J. A., O'Dea, E., Holt, J., Polton, J., Hewitt, H. T., Furner, R., Guihou, K., Brereton, A., Arnold, A.,
1179 Wakelin, S., Castillo Sanchez, J. M., and Mayorga Adame, C. G.: AMM15: a new high-resolution NEMO
1180 configuration for operational simulation of the European north-west shelf, Geosci. Model Dev., 11, 681–
1181 696, <https://doi.org/10.5194/gmd-11-681-2018>, 2018.

1182 [Hernandez, F., Blockley, E., Brassington, G. B., Davidson, F., Divakaran, P., Drévilion, M., Ishizaki, S.,
1183 Garcia-Sotillo, M., Hogan, P. J., Lagemaat, P., Levier, B., Martin, M., Mehra, A., Mooers, C., Ferry, N.,
1184 Ryan, A., Regnier, C., Sellar, A., Smith, G. C., Sofianos, S., Spindler, T., Volpe, G., Wilkin, J., Zaron, E. D.,
1185 and Zhang, A.: Recent progress in performance evaluations and near real-time assessment of
1186 operational ocean products, J. Oper. Oceanogr., 8, 221–238,
1187 <https://doi.org/10.1080/1755876X.2015.1050282>, 2015.](https://doi.org/10.1080/1755876X.2015.1050282)

1188 [Hernandez, F., Smith, G., Baetens, K., Cossarini, G., Garcia-Hermosa, I., Drevillon, M., Maksymczuk, J.,
1189 Melet, A., Regnier, C., and von Schuckmann, K.: Measuring Performances, Skill and Accuracy in
1190 Operational Oceanography: New Challenges and Approaches. In "New Frontiers in Operational
1191 Oceanography", E. Chassignet, A. Pascual, J. Tintore, and J. Verron, Eds., GODAE OceanView, 759-796,
1192 \[doi:10.17125/gov2018.ch29\]\(https://doi.org/10.17125/gov2018.ch29\), 2018.](https://doi.org/10.17125/gov2018.ch29)

1193 Hersbach, H.: Decomposition of the Continuous Ranked Probability Score for Ensemble Prediction
1194 Systems, Wea. Forecasting, 15, 559–570, [https://doi.org/10.1175/1520-
1195 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.

1196 Holt, J., Hyder, P., Ashworth, M., Harle, J., Hewitt, H. T., Liu, H., New, A. L., Pickles, S., Porter, A., Popova,
1197 E. and Allen, J.: Prospects for improving the representation of coastal and shelf seas in global ocean
1198 models, *Geosci. Model Dev.*, *10*, 499-523, 2017.

1199 [Howarth, M. and Pugh, D.: Chapter 4 Observations of Tides Over the Continental Shelf of North-West](#)
1200 [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\),](#)
1201 [1983](#)

1202 [Juza, M., Mourre, B., Lellouche, J. M., Tonani M., and Tintoré, J.: From basin to sub-basin scale](#)
1203 [assessment and intercomparison of numerical simulations in the western Mediterranean Sea](#)
1204 [J. Mar. Syst., *149*, 36-49, <https://doi.org/10.1016/j.jmarsys.2015.04.010>, 2015.](#)

1205 Keil, C. and Craig, G. C.: A Displacement-Based Error Measure Applied in a Regional Ensemble
1206 Forecasting System, *Mon. Wea. Rev.*, *135*, 3248–3259, <https://doi.org/10.1175/MWR3457.1>, 2007.

1207 King, R., While, J., Martin, M. J., Lea, D. J., Lemieux-Dudon, B, Waters, J., O’Dea, E.: Improving the
1208 initialisation of the Met Office operational shelf-seas model. *Ocean Model*, *130*, 1-14, 2018.

1209 [Le Traon PY, Reppucci A, Alvarez Fanjul E, Aouf L, Behrens A, Belmonte M, Bentamy A, Bertino L, Brando](#)
1210 [VE, Kreiner MB, Benkiran M, Carval T, Ciliberti SA, Claustre H, Clementi E, Coppini G, Cossarini G, De](#)
1211 [Alfonso Alonso-Muñoyerro M, Delamarche A, Dibarboure G, Dinessen F, Dreviron M, Drillet Y, Faugere](#)
1212 [Y, Fernández V, Fleming A, Garcia-Hermosa MI, Sotillo MG, Garric G, Gasparin F, Giordan C, Gehlen M,](#)
1213 [Gregoire ML, Guinehut S, Hamon M, Harris C, Hernandez F, Hinkler JB, Hoyer J, Karvonen J, Kay S, King R,](#)
1214 [Lavergne T, Lemieux-Dudon B, Lima L, Mao C, Martin MJ, Masina S, Melet A, Buongiorno Nardelli B,](#)
1215 [Nolan G, Pascual A, Pistoia J, Palazov A, Piolle JF, Pujol MI, Pequignet AC, Peneva E, Pérez Gómez B, Petit](#)
1216 [de la Villeon L, Pinardi N, Pisano A, Pouliquen S, Reid R, Remy E, Santoleri R, Siddorn J, She J, Staneva J,](#)
1217 [Stoffelen A, Tonani M, Vandenbulcke L, von Schuckmann K, Volpe G, Wettre C and Zacharioudaki A:](#)
1218 [From Observation to Information and Users: The Copernicus Marine Service Perspective. *Front. Mar. Sci.*](#)
1219 [6:234. doi: 10.3389/fmars.2019.00234, 2019.](#)

1220 [Lorente, P., Sotillo, M., Amo-Baladrón, A., Aznar, R., Levier, B., Aouf, L., Dabrowski, T., Pascual, Á.,](#)
1221 [Reffray, G., Dalphinnet, A., Toledano Lozano, C., Rainaud, R., and Alvarez Fanjul, E. : The NARVAL Software](#)
1222 [Toolbox in Support of Ocean Models Skill Assessment at Regional and Coastal Scales.](#)
1223 http://doi.org/10.1007/978-3-030-22747-0_25 2019a.

1224 [Lorente, P., García-Sotillo, M., Amo-Baladrón, A., Aznar, R., Levier, B., Sánchez-Garrido, J. C.,](#)
1225 [Sammartino, S., de Pascual-Collar, Á., Reffray, G., Toledano, C., and Álvarez-Fanjul, E.: Skill assessment of](#)
1226 [global, regional, and coastal circulation forecast models: evaluating the benefits of dynamical](#)
1227 [downscaling in IBI \(Iberia–Biscay–Ireland\) surface waters, *Ocean Sci.*, 15, 967–996,](#)
1228 <https://doi.org/10.5194/os-15-967-2019>, 2019b.

1229 Madec, G. and the NEMO team: NEMO ocean engine. Note du Pôle de modélisation, Institut Pierre-
1230 Simon Laplace (IPSL), France, No 27 ISSN No 1288-1619, 2016.

1231 [Mason, E., Ruiz, S., Bourdalle-Badie, R., Reffray, G., García-Sotillo, M., and Pascual, A.: New insight into](#)
1232 [3-D mesoscale eddy properties from CMEMS operational models in the western Mediterranean, *Ocean*](#)
1233 [Sci.](#), 15, 1111–1131, <https://doi.org/10.5194/os-15-1111-2019>, 2019.

1234 Mass, C. F., Ovens, D., Westrick, K., and Colle, B. A.: DOES INCREASING HORIZONTAL RESOLUTION
1235 PRODUCE MORE SKILLFUL FORECASTS?, *Bull. Amer. Meteor. Soc.*, 83, 407–430,
1236 [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.

1237 Mirouze, I., Blockley, E. W., Lea, D. J., Martin, M. J. and Bell, M. J.: A multiple length scale correlation
1238 operator for ocean data assimilation, *Tellus A*, 68, 29744, <https://doi.org/10.3402/tellusa.v68.29744>,
1239 2016.

1240 Mittermaier, M., Roberts, N., and Thompson, S. A.: A long-term assessment of precipitation forecast skill
1241 using the Fractions Skill Score, *Met. Apps*, 20, 176-186, <https://doi.org/10.1002/met.296>, 2013.

1242 Mittermaier, M. P.: A Strategy for Verifying Near-Convection-Resolving Model Forecasts at Observing
1243 Sites, *Wea. Forecasting*, 29, 185–204, <https://doi.org/10.1175/WAF-D-12-00075.1>, 2014.

1244 Mittermaier, M. P., and Csima, G.: Ensemble versus Deterministic Performance at the Kilometer Scale,
1245 *Wea. Forecasting*, 32, 1697–1709, <https://doi.org/10.1175/WAF-D-16-0164.1>, 2017.

1246 Mogensen, K, Balmaseda, M. A., Weaver, A.: The NEMOVAR ocean data assimilation system as
1247 implemented in the ECMWF ocean analysis for System 4. European Centre for Medium-Range Weather
1248 Forecasts, 2012.

1249 [Mourre B., E. Aguiar, M. Juza, J. Hernandez-Lasheras, E. Reyes, E. Heslop, R. Escudier, E. Cutolo, S. Ruiz,](#)
1250 [E. Mason, A. Pascual and J. Tintoré: Assessment of high-resolution regional ocean prediction systems](#)
1251 [using multi-platform observations: illustrations in the Western Mediterranean Sea. In “New Frontiers in](#)
1252 [Operational Oceanography”, E. Chassignet, A. Pascual, J. Tintoré and J. Verron, Eds, GODAE Ocean View,](#)
1253 [663-694, doi: 10.17125/gov2018.ch24, 2018.](#)

1254 O'Dea, E. J., Arnold, A. K., Edwards, K. P., Furner, R., Hyder, P., Martin, M. J., Siddorn, J. R., Storkey, D.,
1255 While, J., Holt, J. T., and Liu, H.: An operational ocean forecast system incorporating NEMO and SST data
1256 assimilation for the tidally driven European North-West shelf, *J. Oper. Oceanogr.*, 5, 3–17,
1257 <https://doi.org/10.1080/1755876X.2012.11020128>, 2012.

1258 O'Dea, E., Furner, R., Wakelin, S., Siddorn, J., While, J., Sykes, P., King, R., Holt, J., and Hewitt, H.: The
1259 CO5 configuration of the 7 km Atlantic Margin Model: large-scale biases and sensitivity to forcing,
1260 physics options and vertical resolution, *Geosci. Model Dev.*, 10, 2947–2969,
1261 <https://doi.org/10.5194/gmd-10-2947-2017>, 2017.

1262 Rossa A., Nurmi P., Ebert E.: Overview of methods for the verification of quantitative precipitation
1263 forecasts, in: *Precipitation: Advances in Measurement, Estimation and Prediction*, edited by:
1264 Michaelides, S., Springer, Berlin, Heidelberg, 419–452, https://doi.org/10.1007/978-3-540-77655-0_16,
1265 2008.

1266 Tonani, M., Sykes, P., King, R. R., McConnell, N., Péquignat, A.-C., O'Dea, E., Graham, J. A., Polton, J., and
1267 Siddorn, J.: The impact of a new high-resolution ocean model on the Met Office North-West European
1268 Shelf forecasting system, *Ocean Sci.*, 15, 1133–1158, <https://doi.org/10.5194/os-15-1133-2019>, 2019.

1269
1270 World Meteorological Organisation: Guide to Meteorological Instruments and Methods of Observation
1271 (WMO-No. 8, the CIMO Guide) –available at
1272 https://library.wmo.int/opac/doc_num.php?explnum_id=4147, 2017.

1273 9. Author contributions

1274 All authors contributed to the introduction, data and methods, and conclusions. RC, JM and MM
1275 contributed to the scientific evaluation and analysis of the results. RC and JM designed and ran

1276 the model assessments. CP supported the assessments through the provision and reformatting
1277 of the data used. MT provided detail on the model configurations used.

1278

1279 **10. Competing interests**

1280 The authors declare that they have no conflict of interest.

1281

1282 **11. Acknowledgements**

1283 This study has been conducted using E.U. Copernicus Marine Service Information.

1284 This work has been carried out as part of the Copernicus Marine Environment Monitoring Service
1285 (CMEMS) **HiVE** project. CMEMS is implemented by Mercator Ocean International in the
1286 framework of a delegation agreement with the European Union.

1287 Model Evaluation Tools (MET) was developed at the National Center for Atmospheric Research
1288 (NCAR) through grants from the National Science Foundation (NSF), the National Oceanic and
1289 Atmospheric Administration (NOAA), The United States Air Force (USAF), and the United States
1290 Department of Energy (DOE). NCAR is sponsored by the United States National Science
1291 Foundation.

1292