

1 Model performance

The following section intends to explore the performance and capabilities of the deep learning system. Where the previous section 1.3 assessed the intra-training model performance, the current section will compare a benchmark deep learning model against baselines and physical models. The physical models have been previously described in section 1, and the baselines (although previously mentioned and to some extent utilized) will be derived in the following subsection. This section will first assess model performance against persistence. Afterwards, the deep learning system will be compared against other physical models. Setup and considerations will be described as they become relevant.

1.1 Baselines

Two types of baselines are considered, persistence and a linear trend. A persistence forecast is constant in time. Regardless of the forecast lead time the initial values for all grid cells are kept constant. Moreover, the autocorrelation of sea ice concentration from the sea ice charts was shown in section 1.2.1 to be high for short lead times. To determine if a forecast has predictive skill, the forecast has to achieve a lower NIIEE than persistence, which is a similar approach as employed in Zampieri et al. (2019). We believe that using this threshold as the definition of a skillful forecast preserves the intent of validating the sea ice forecast in a manner relevant for maritime end users (Melsom et al., 2019; Veland et al., 2021).

The second baseline uses the linear trend, as described in section 1.2.3 and used as predictor for the deep learning system 1.1.3. However, the computed linear trend will be applied pixelwise to advance the initial state forward in time to a given lead time. As the linear trend is computed from OSI SAF ssmis observations, it will consequently be applied to the same dataset. For clarity, the linear trend forecast is computed on the 1km AROME Arctic grid, and the computed values are clipped to match the valid value range, i.e. $\text{values} < 0 \rightarrow \text{values} = 0 \wedge \text{values} > 100 \rightarrow \text{values} = 100$.

1.2 Verifying performance against persistence

For this section, a model representing a benchmark with a depth of 256 channels in the final feature map, learning rate = 0.001 and all predictor variables have been used. Only the core training dataset was used for training(2019 and 2020).

The seasonal distribution of average ice edge displacement for all sea ice categories found in the sea ice charts are shown for the deep learning system and persistence are displayed

in figure 1. Figure 1 demonstrates the predictive performance for the deep learning system measured at each resolved contour. In figure 1 b), c), d) and e), the deep learning system achieves a lower median 25-th and 75-th percentile than persistence.

Figures 2 and 3 shows the model confidence as an annual mean for all output contours (figure 2) and the ($> 10\%$) contour distributed seasonally (figure 3). The confidence values shown are output pixel values after the sigmoid (equation 7), such that values closer than 1 are pixels that the model is more confident to belong in the outputted contour. Likewise, values closer to 0 are confident not to belong to the targeted contour.

Figure 2 show that all cumulative contours, except $= 100\%$, have a similar confidence pattern similar to the seasonal cycle of sea ice concentration. Which is expected since for all contours at all dates the leftmost sea ice concentration is always present, whereas the less confident areas east of Svalbard exert spatial variability through mobility. However, it is noted that the $= 100\%$ contour (figure 2 (f)) show a lower overall confidence level, as well as only being restricted to the land structures present in the scene.

The seasonal confidence cycle for the $= 100\%$ contour is shown in figure 3. It is seen that the spatial distribution of confidence tend cover the land-covered pixels for all seasons, although with varying levels of confidence.

The monthly mean sea ice edge length is shown in figure 4. From figure 4, the predicted ice edge follow a similar seasonal pattern to the ice edge length from the target ice charts. Each monthly mean predicted sea ice edge length is biased towards shorter lengths, with the annual mean bias for 2022 being -2146km.

Moreover, the monthly distribution of the different sea ice categories is shown in figure 5. The figure show that the deep learning resolve the area of each contour with a similar scale and variability as the target sea ice charts.

1.3 Inter-product comparison

This section covers results regarding the multi-product comparison. First, the preparation of samples as well as setup of the comparison environment is described. The physical models considered for this comparison are neXtSIM (Williams et al., 2021) presented in section 1.3.2 and Barents-2.5 (Röhrs et al., 2022) presented in section 1.3.3, whereas the considered baselines are persistence and the linear sea ice concentration trend described in section 1.1. Two different products are used as targets. The first product are the sea ice charts, which will be utilized similarly as when comparing only against persistence in section 1.2. Secondly, the independent AMSR2 observations produced by Spreen et al. (2008) are also utilized as ground truth targets.

Seasonal distribution of average ice edge displacement

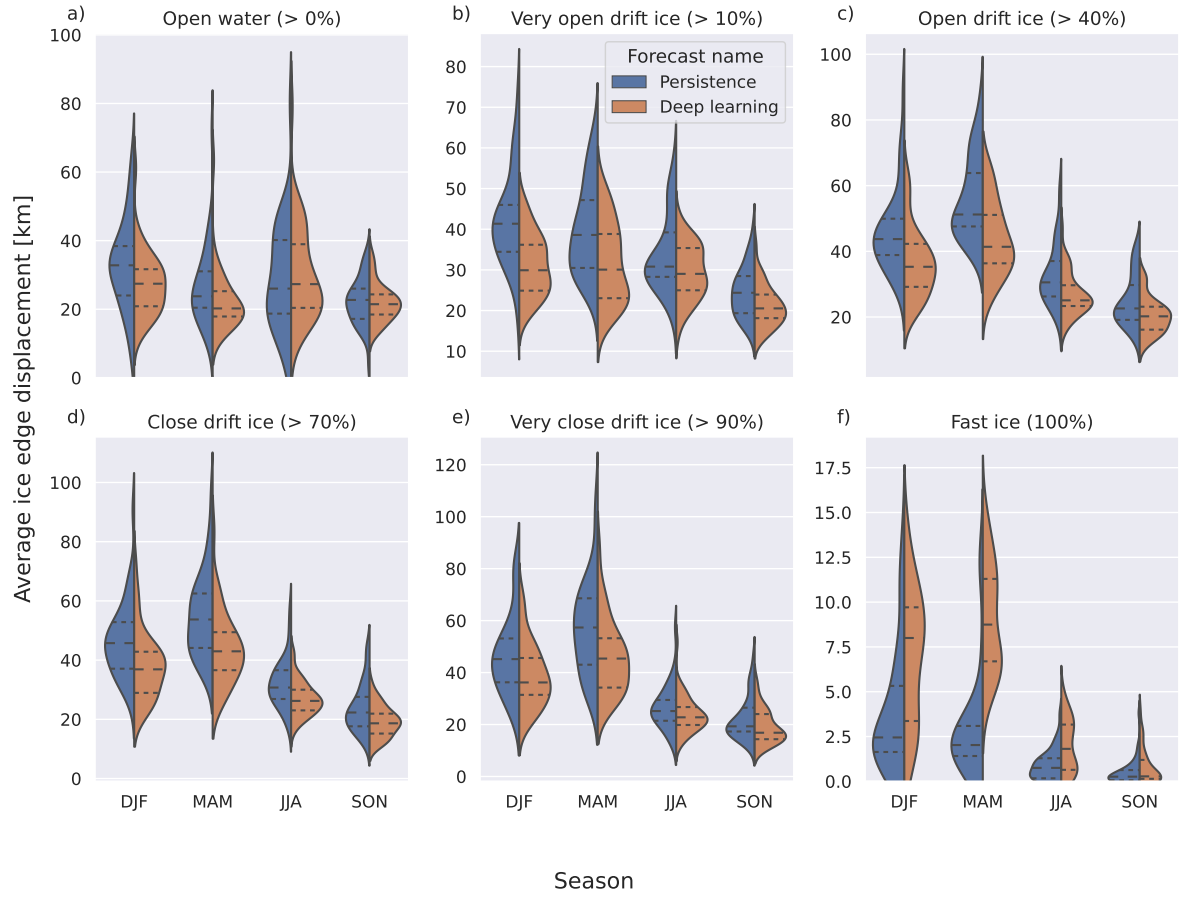


Figure 1: Seasonal distribution of the mean ice edge displacement (Normalized IIEE) for the different sea ice chart categories in the form of cumulative contours. The related sea ice concentration range for each contour is also included. The lower and upper dashed line denote the interquartile range, with the middlemost dashed line showing the distribution median.

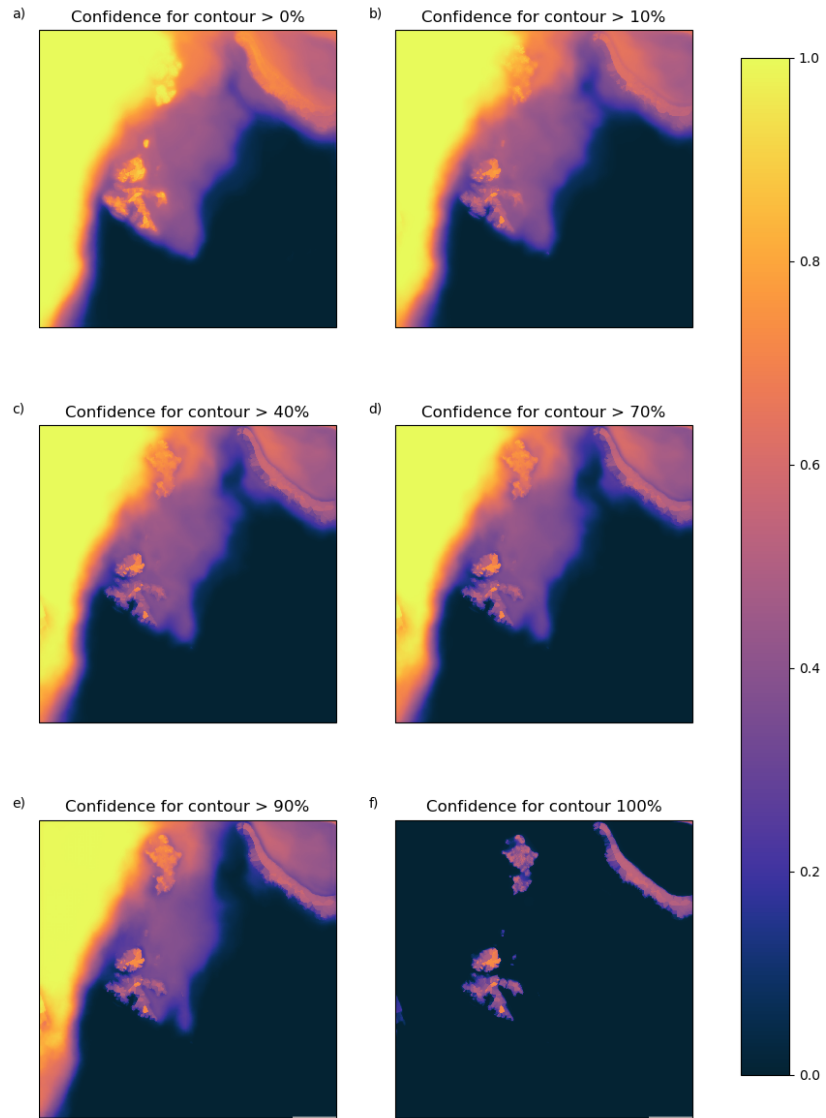


Figure 2: Mean annual probabilities for the different cumulative contours outputted by the model (the class ice free open water is not shown).

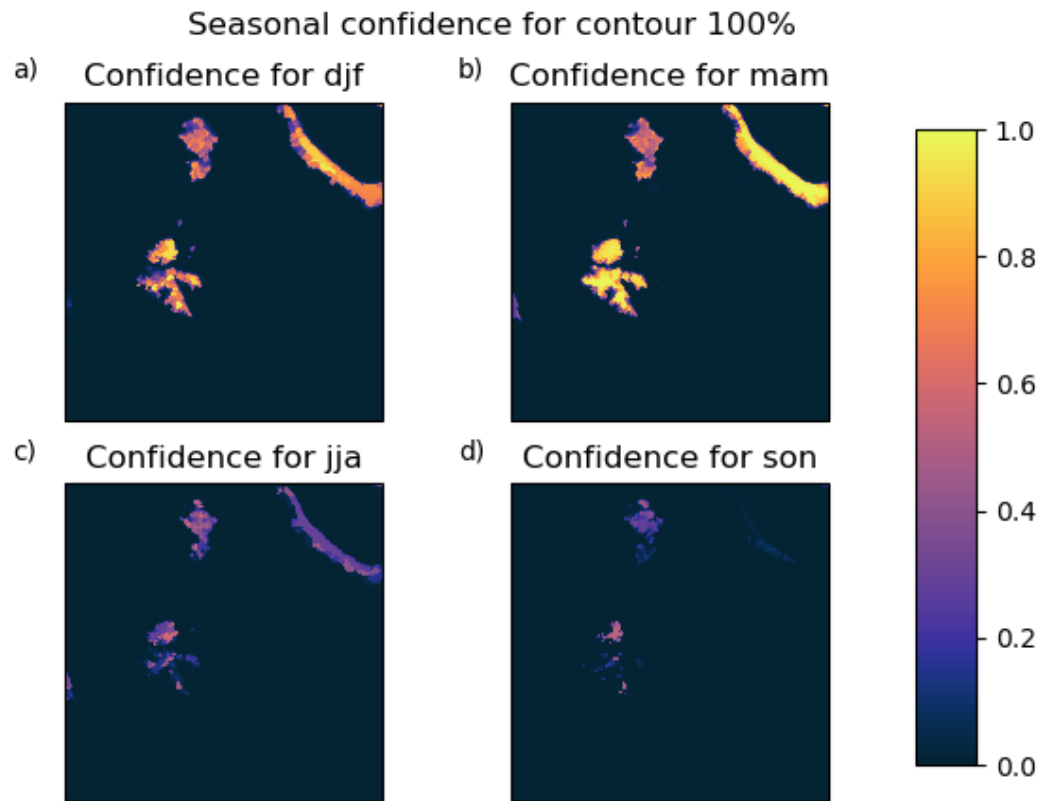


Figure 3: Mean seasonal confidence for the ($= 100\%$) cumulative contour.

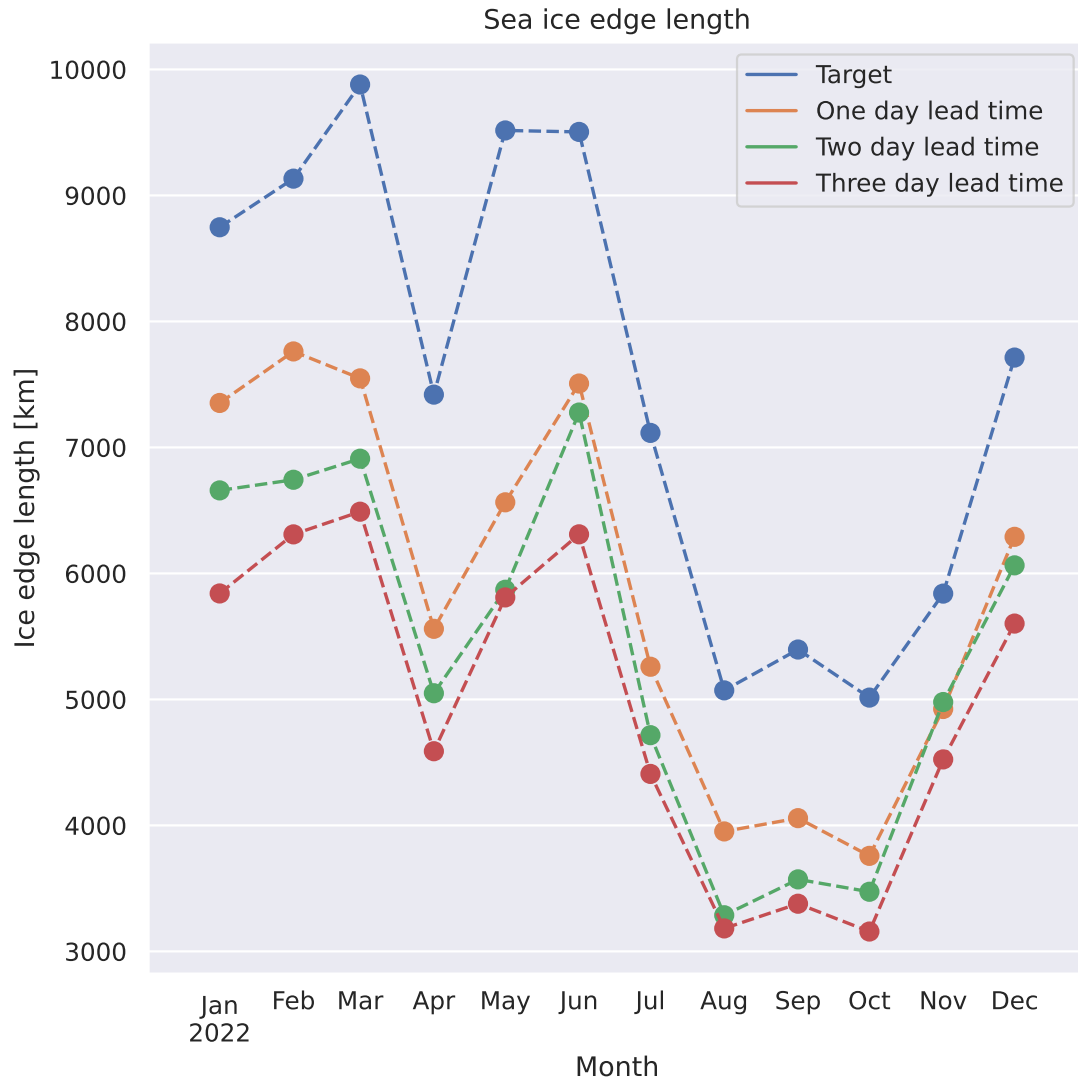


Figure 4: Mean monthly sea ice edge length for the entire 2022 test dataset. The ice edge is defined from a 10% threshold, which results in the ($> 10\%$) contour being used to define the ice edge. Each entry in the defined sea ice edge are on a 1km resolution. Each deep learning marker is annotated with the mean monthly bias with respect to the target sea ice edge length.

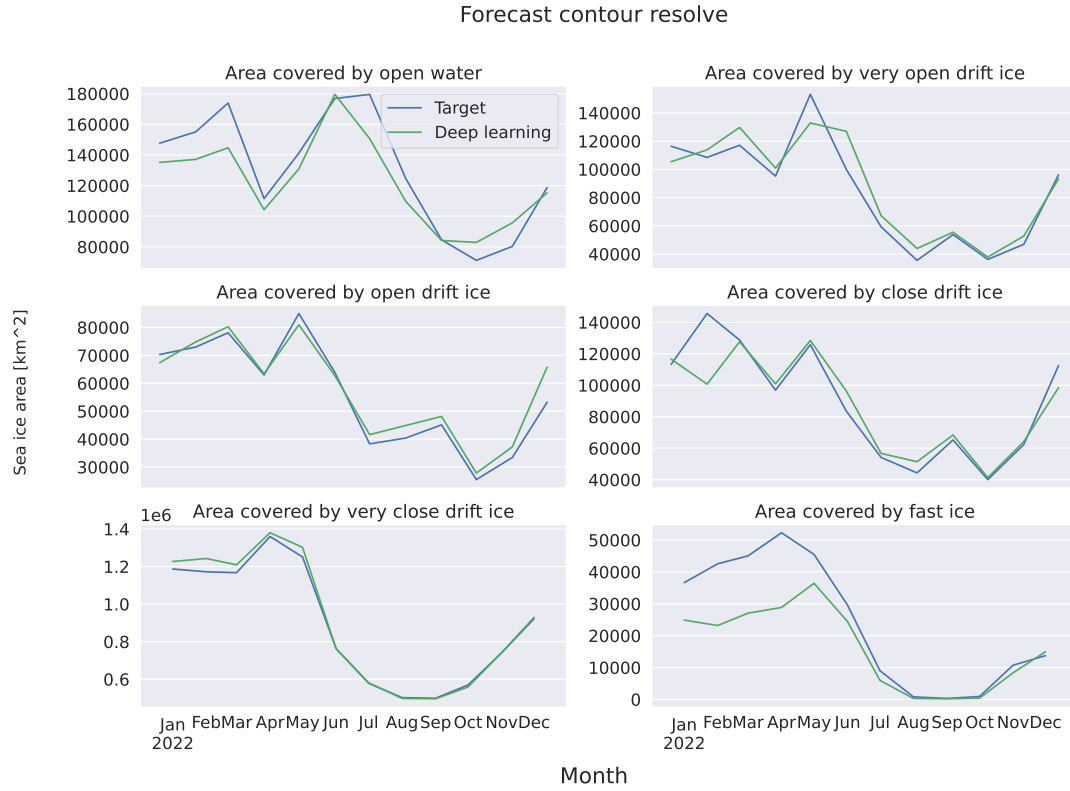


Figure 5: Mean monthly sea ice category distribution for the model and the target sea ice charts for the 2022 test dataset. Each contour is represented by the sea ice area, which is computed from the sum of pixels in each contour times their spatial extent.

When comparing against multiple products, the coarsest resolution model is used as a common spatial resolution. Also, the projection of the coarsest resolution is used for all products, such that other products have to be interpolated onto the grid of the coarsest resolution model, which is done using nearest neighbor interpolation. As both baselines have a daily forecast frequency, comparing either with a deep learning prediction involves identifying the forecast with similar bulletin- and valid date, i.e. initialized at the same day and targeting the same lead time. When utilizing the sea ice charts as the ground truth, the spatial resolution of neXtSIM (3km) is the coarsest, and thus all products are interpolated onto the same resolution.

Comparing against the two physical models requires a consideration of the hourly forecast frequency (Williams et al., 2021; Röhrs et al., 2022) of both model. First, given a published sea ice chart, the comparable physical model is initialized the following day at 00:00 UTC. Furthermore, a daily mean is computed from the 24 steps forward in time taken by the physical model when it covers the valid date of the deep learning forecast. Even though the sea ice charts only convey information about the sea ice concentration up until their publication time, the operational product is considered a reference for the entirety of the publication date. Moreover, to reduce introducing a bias towards the time of day to the physical forecasts as well as limiting the spatial variability induced by the lack of a temporal mean, reducing the physical forecasts to daily averages is considered a more comparable approach than e.g. selecting a single hour (15:00 UTC) from the forecasts.

Since the AMSR2 observations are supplied on a 6.25 km spatial resolution (Spreen et al., 2008), when AMSR2 is used as the ground truth all data is interpolated to match the resolution of AMSR2. Although the AMSR2 data have a substantially coarser spatial resolution compared to the sea ice charts or the deep learning system, the data makes it possible to assess the generalizability of the deep learning performance when targeting an unseen ground truth.

From this setup, the mean of the first 24 hours of a forecast from a physical model is compared against a deep learning prediction with one day lead time, the mean between 24 and 48 hours are compared against a deep learning prediction with two day lead and the mean of the third predicted day is compared against a deep learning prediction with three day lead time. Figure 6 summarizes the process. Note that Barents-2.5 only have a 66 hour lead time (Röhrs et al., 2022), thus the mean between $t = 48$ and $t = 66$ is computed when comparing against a three day lead time prediction.

It is noted that when comparing against multiple forecast products as described in figure 6, only the common dates shared between all products are used. With the current setup, where neXtSIM, Persistence, Deep learning, OSI SAF trend and Barents-2.5 are considered, the test dataset is reduced from 196 to 171 samples, 147 to 130 samples and 142

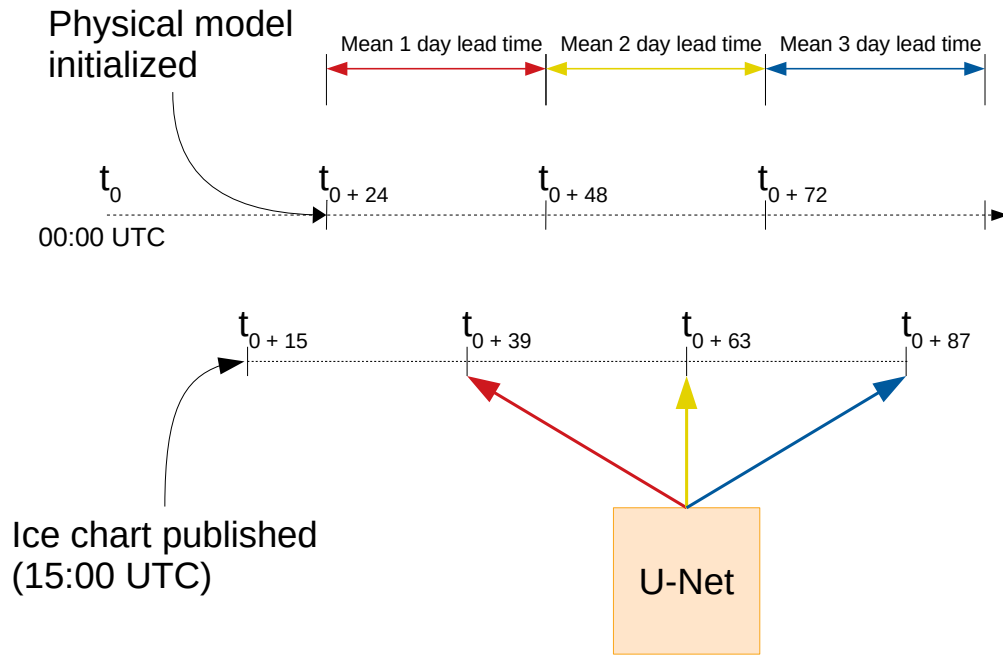


Figure 6: Overview describing how a physical model with an hourly frequency is compared against a deep learning forecast. Timestamps are hourly, and relative to 00:00 UTC the day a sea ice chart is published. The physical model is initialized the following day. Colors are used to denote lead time comparability, with red = 1, yellow = 2 and green = 3 day lead time.

to 125 samples for 1, 2, and 3 day lead time respectively. Moreover, Barents-2.5 is only considered starting with the month of June, to comply with the spin up time of its data assimilation system (Röhrs et al., 2022).

Figure 7 shows the seasonal distribution of NIIIE for the different forecast systems and benchmarks, following the setup described in figure 6. By inspecting figure 7, it can be seen that only the products based on the sea ice chart are able to achieve consistently low NIIIE for the $> 0\%$ contour. Furthermore, for the $\geq (10, 40, 70, 90)\%$ contours, the deep learning system achieves the lowest median and mean values compared to all the other products. It can also be seen that neXtSIM tend to increase its mean and median as well as spread for increasing contours, with a similar although not as consistent pattern for Barents-2.5. Moreover, the OSI SAF trend typically have the highest (visible) outliers. Finally, no product is able to achieve a lower mean or median NIIIE compared to persistence when inspecting the 100% (fast ice) contour.

The fraction of days where Deep learning achieves lower NIIIE compared to each considered product is shown in Figure 8. The figure shows that the deep learning system consistently achieves a $\geq 50\%$ success rate compared to all products, except for Persistence 1 day lead time in July, August and September as well as Barents-2.5 2 day lead time in November and December. When compared to neXtSIM at 1 day lead time (figure 8 (a)), the Deep learning system achieves a lower NIIIE at all considered dates in the test data. However, it can also be seen that a lower amount of days with lower NIIIE than neXtSIM are achieved as the lead time increases. The same pattern may also be seen in the Barents data as the mean fraction of days with lower NIIIE for the Deep learning system also decrease with lead time, although Barents is only able to achieve lower NIIIE more than 50% of the dates for a 2 day lead time as previously noted. With respect to persistence, the Deep learning forecasts seem to achieve a higher fraction of days with lower NIIIE as lead time increases, although there is no trend for the individual months. At the ($\geq 10\%$) contour, the OSI SAF trend is consistently beat by the deep learning system during Winter and Spring, with less consistency observed during the Summer and Autumn seasons.

The spatial distribution of product error is shown in Figure 9. From the figure, it can be seen that both sea ice chart products have lower bias than the three other products, as well as only exerting biases in the MIZ. Moreover, it can be seen from the top row in figure 9 the neXtSIM data have a negative bias along the sea ice edge, which is prominent during Winter and Spring. Moreover, the OSI SAF trend seem to have a strong negative bias along a wide sea ice edge. Finally, Barents seem to have a positive bias around Svalbard in the Summer, with a less prominent overall bias during the Fall.

The seasonal NIIIE distributions shown in figure 10 is created similarly as figure 7, but with AMSR2 as the ground truth data, which also implies that all data have been interpo-

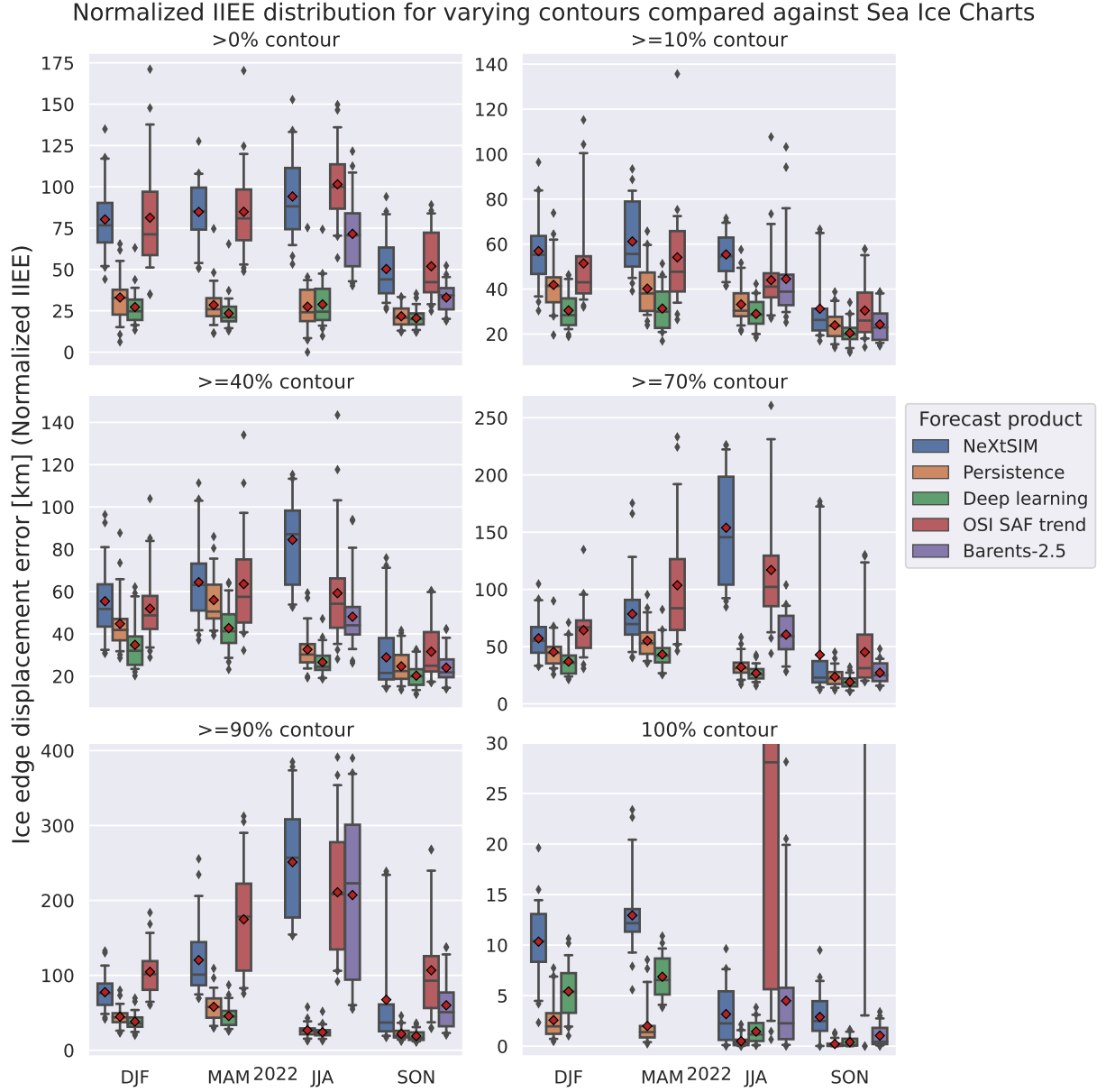


Figure 7: Model intercomparison with a two day lead time. The boxes are constructed from seasonally distributed NIIEE values computed from the test dataset (2022). The sea ice charts are considered as targets. Each box cover the interquartile range (25th - 75th percentile), with whiskers covering the 5th and 95th percentile. The line in each box is the median, and the red diamond is the mean. The IIEE is normalized according to the climatological sea ice edge at the forecast valid date. The extent of the y axis is limited in such a way that the distributions are easily readable, at the expense of some outliers not being visible. The OSI SAF trend is computed from the past 7 days.

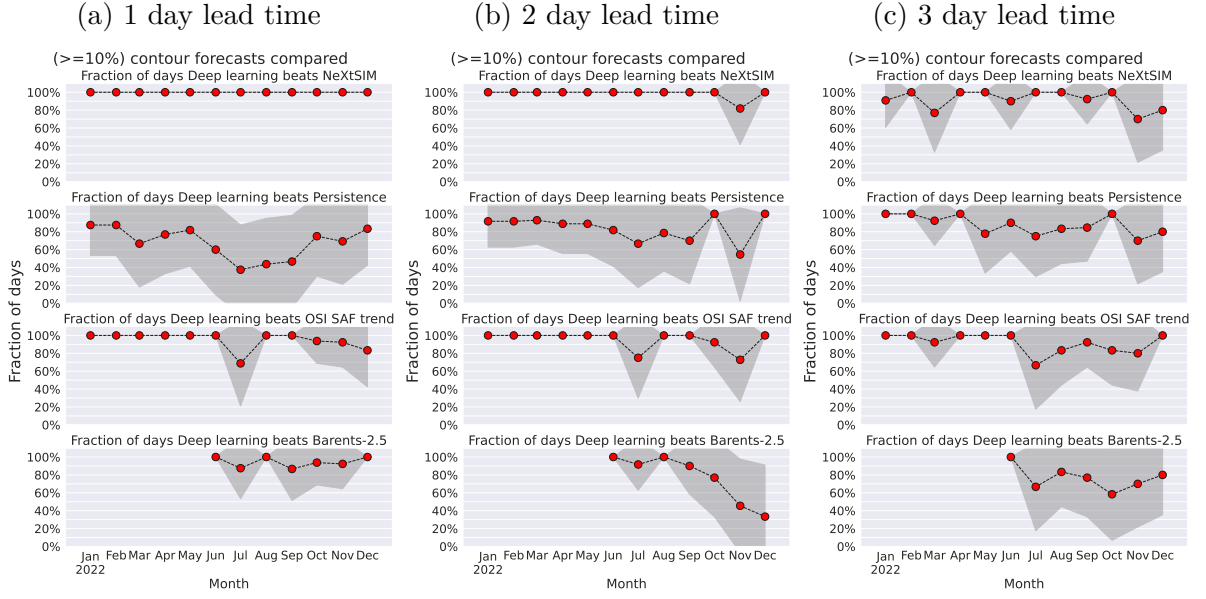


Figure 8: Fraction of days where the Deep learning forecast achieves a lower NIIIE than the compared product, distributed monthly for all lead times. Only the ($\geq 10\%$) contour has been considered, due to the relevance of the contour with respect to the definition of the sea ice edge and its application to operational end users. Gray contours denote the uncertainty (standard deviation) for each month. The sea ice chart has been used as ground truth target when computing the IIEE, and the score has been normalized according to the climatological sea ice edge.

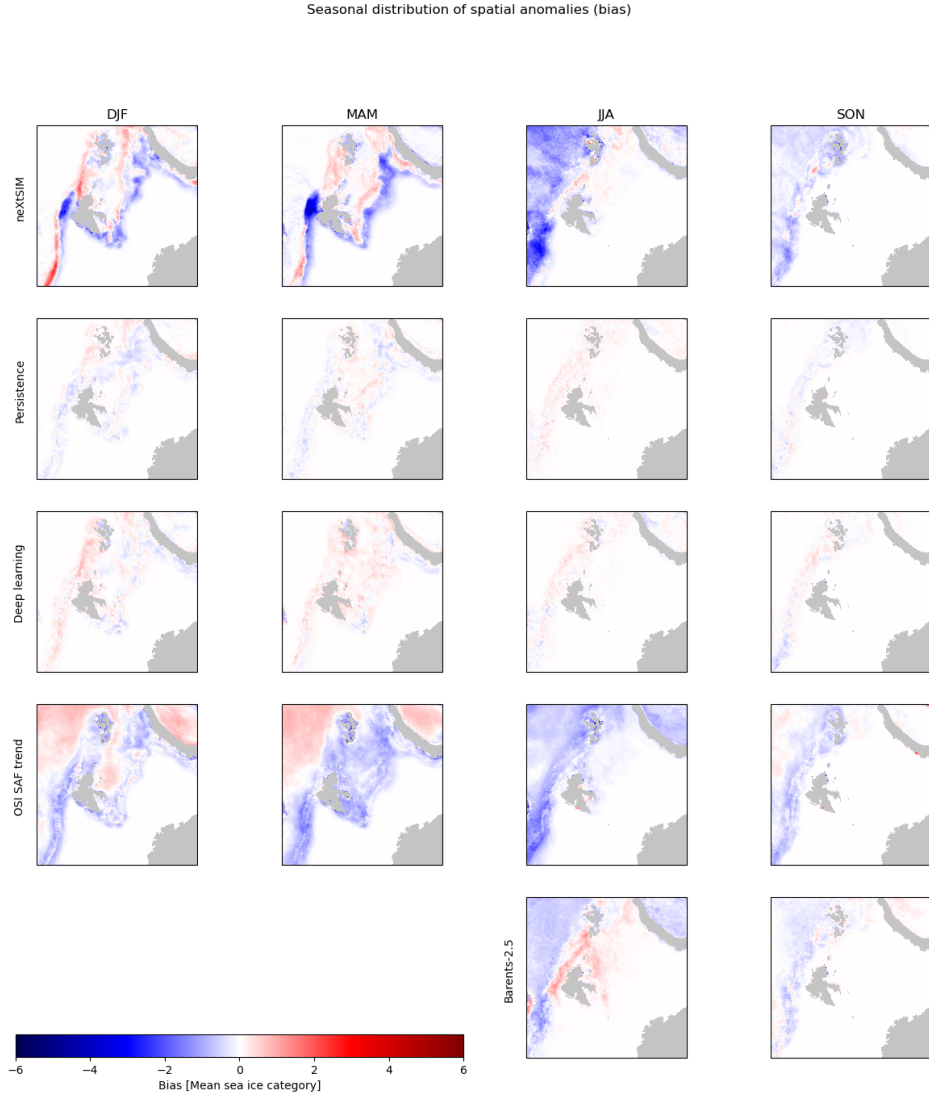


Figure 9: Spatial distribution of the mean seasonal error in predicted sea ice categories between the compared products. The data is interpolated onto the neXtSIM grid, and the test dataset is considered.

lated onto the 6.25km AMSR2 grid (Spreen et al., 2008). Contrary to what was observed in figure 7, the sea ice chart based products in figure 10 exert significantly higher NIIIE at the $> 0\%$ contour. However, Barents-2.5 also exert a similar increased NIIIE as the sea ice chart based products at the same contour. Moreover, the sea ice chart based models are within the interquartile range of neXtSIM and OSI SAF trend starting at the ($\geq 10\%$) contour. At the 0 and 10% contours, the OSI SAF trend exerts the lowest mean and median NIIIE for all months except SON where neXtSIM achieves the lowest median and mean. However, starting at the ($\geq 40\%$) the deep learning system has the lowest median and mean NIIIE, which lasts until the 100% contour where performance is comparable between all products except for the OSI SAF trend during DJF and MAM.

Following the result seen in the upper leftmost distribution in Figure 10, Figure 11 shows a comparable figure but with a deep learning model which does not predict the 10% and 100% contours as described in section 1.3.4. By inspecting the $>0\%$ contour, it can be seen that the deep learning system achieves significantly lower NIIIE than persistence, as well as the deep learning system in figure 10. Otherwise for the other contours, the performance of the deep learning system is comparable to the deep learning system in figure 10.

The boxplots in Figure 12 computes the $>0\%$ contour NIIIE against AMSR2 with the model used in Figure 10 but with the predicted $>0\%$ contour removed. The distribution seen in the figure resembles that in Figure 11, with the Deep learning forecasts performing significantly better than persistence.

Kanskje dette også skal i appendix?

References

- Melsom, A., Palerme, C., and Müller, M.: Validation metrics for ice edge position forecasts, *Ocean Science*, 15, 615–630, <https://doi.org/10.5194/os-15-615-2019>, 2019.
- Röhrs, J., Gusdal, Y., Rikardsen, E., Moro, M. D., Brændshøi, J., Kristensen, N. M., Fritzner, S., Wang, K., Sperrevik, A. K., Idžanović, M., Lavergne, T., Debernard, J., and Christensen, K. H.: "in prep for GMD" An operational data-assimilative coupled ocean and sea ice ensembleprediction model for the Barents Sea and Svalbard, p. 20, 2022.
- Spreen, G., Kaleschke, L., and Heygster, G.: Sea ice remote sensing using AMSR-E 89-GHz channels, *Journal of Geophysical Research*, 113, <https://doi.org/10.1029/2005jc003384>, 2008.
- Veland, S., Wagner, P., Bailey, D., Everett, A., Goldstein, M., Hermann, R., Hjort-Larsen, T., Hovelsrud, G., Hughes, N., Kjøl, A., Li, X., Lynch, A., Müller, M., Olsen, J., Palerme, C., Pedersen, J. L., Rinaldo, ., Stephenson, S., and Storelvmo, T.: Knowledge

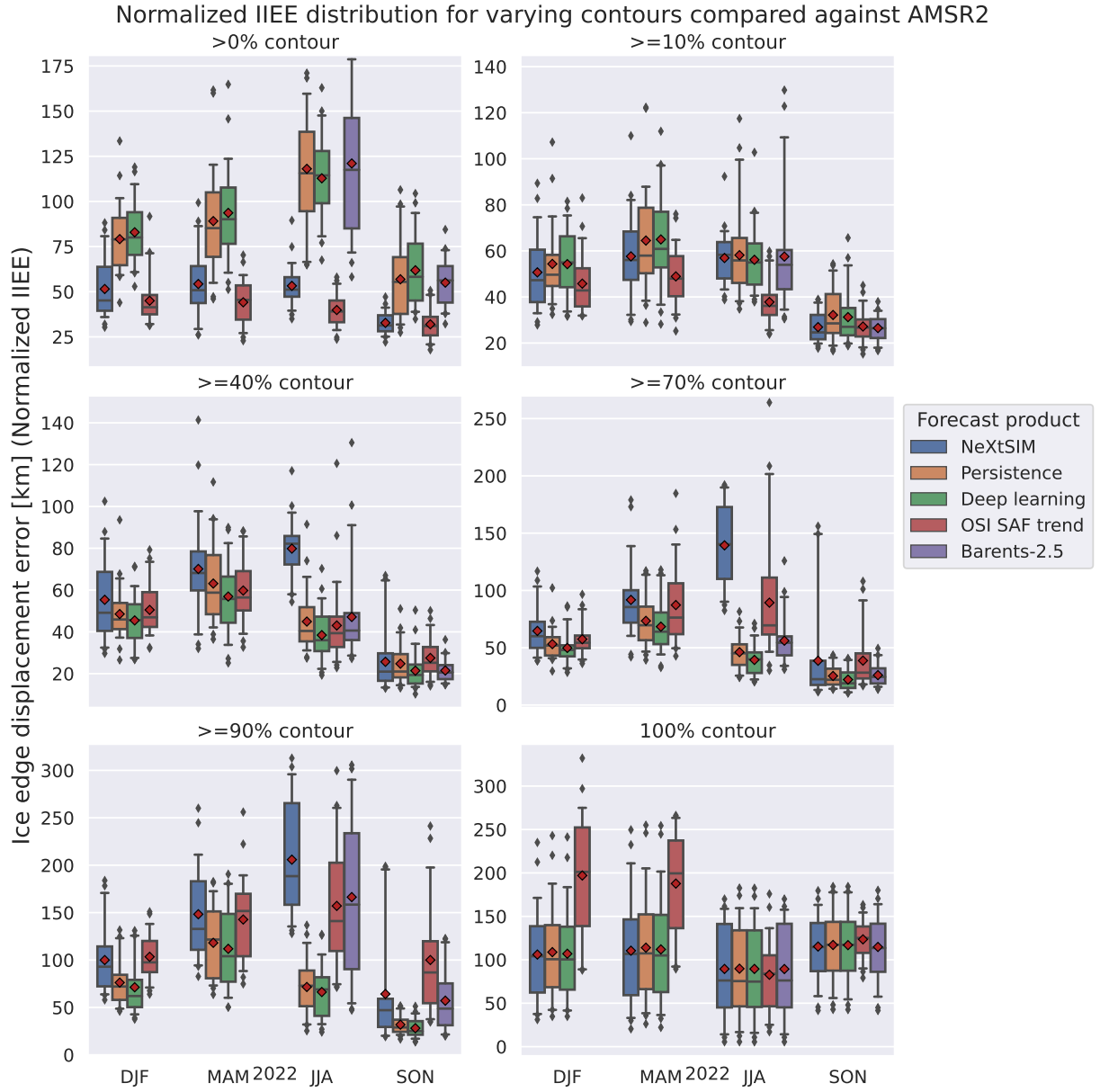


Figure 10: Same as figure 7, but with AMSR2 sea ice concentration as the target ground truth data.

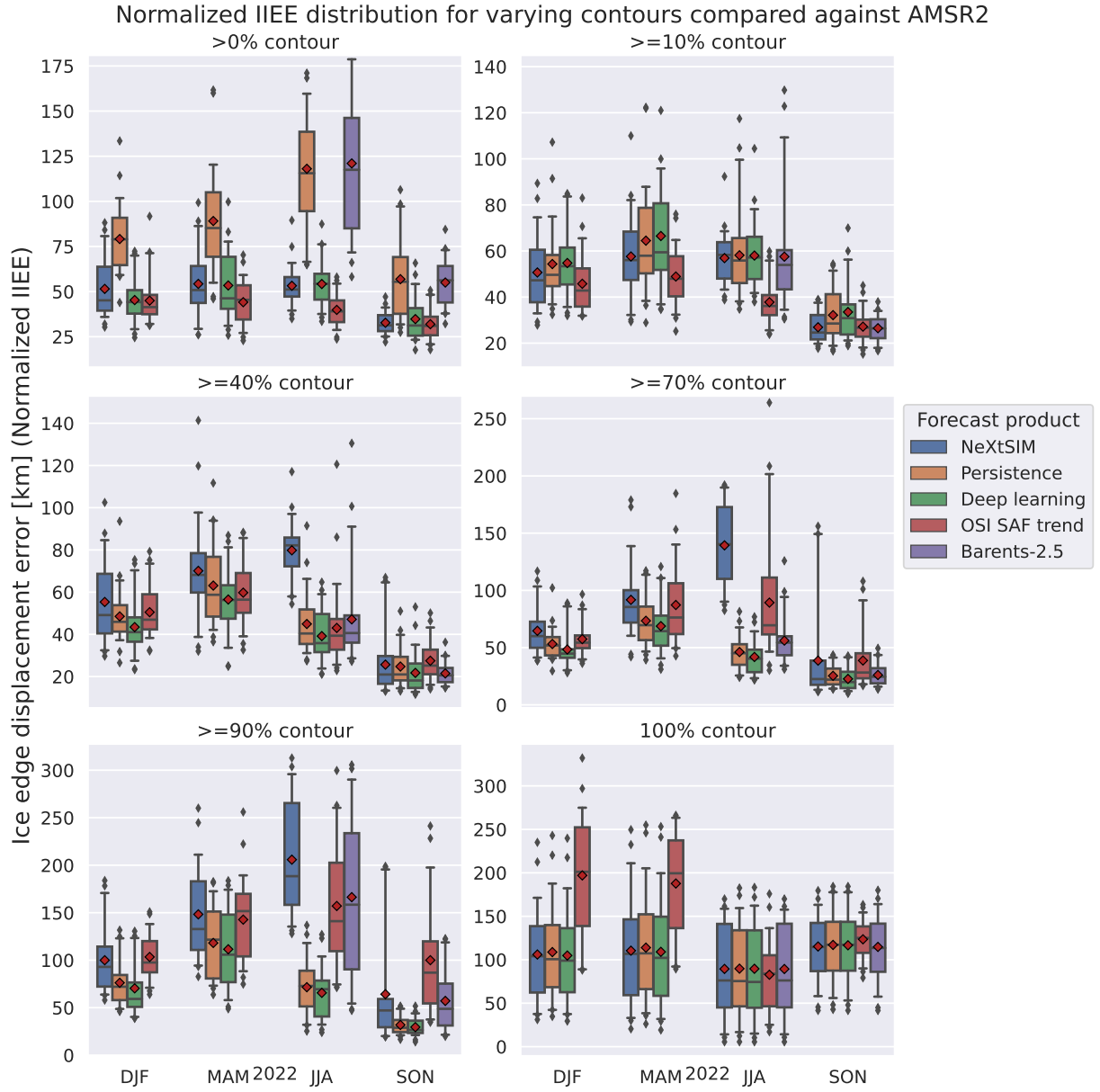


Figure 11: Same as figure 10, but the deep learning system used has reduced output classes.

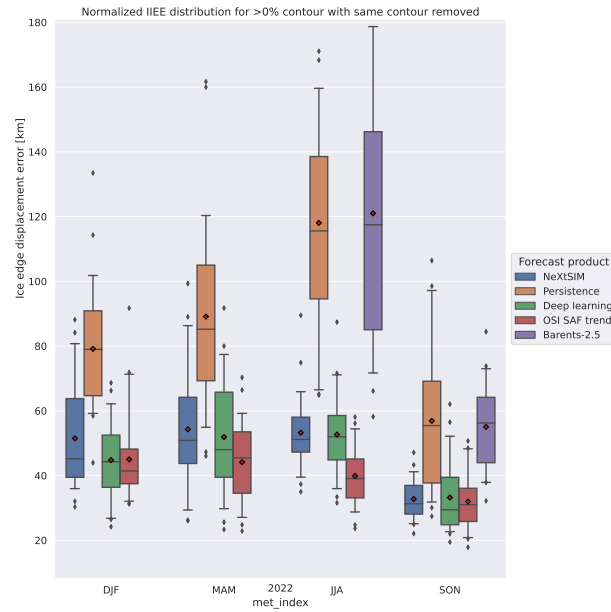


Figure 12: NIIEE for the >0% contour with the model from Figure 10, but with the values in the >0% contour set to category 0 (ice free open water)

needs in sea ice forecasting for navigation in Svalbard and the High Arctic, Tech. Rep. NF-rapport 4/2021, Svalbard Strategic Grant, Svalbard Science Forum, 2021.

Williams, T., Korosov, A., Rampal, P., and Ólason, E.: Presentation and evaluation of the Arctic sea ice forecasting system neXtSIM-F, *The Cryosphere*, 15, 3207–3227, <https://doi.org/10.5194/tc-15-3207-2021>, 2021.

Zampieri, L., Goessling, H. F., and Jung, T.: Predictability of Antarctic Sea Ice Edge on Subseasonal Time Scales, *Geophysical Research Letters*, 46, 9719–9727, <https://doi.org/10.1029/2019gl084096>, 2019.