

Recognition of polar lows in Sentinel-1 SAR images with deep learning

Jakob Grahn*, Filippo Maria Bianchi

Abstract—In this paper, we explore the possibility of detecting polar lows in C-band SAR images by means of deep learning. Specifically, we introduce a novel dataset consisting of Sentinel-1 images divided into two classes, representing the *presence* and *absence* of a maritime mesocyclone, respectively. The dataset is constructed using the ERA5 dataset as baseline and it consists of 2004 annotated images. To our knowledge, this is the first dataset of its kind to be publicly released. The dataset is used to train a deep learning model to classify the labeled images. Evaluated on an independent test set, the model yields an F-1 score of 0.95, indicating that polar lows can be consistently detected from SAR images. Interpretability techniques applied to the deep learning model reveal that atmospheric fronts and cyclonic eyes are key features in the classification. Moreover, experimental results show that the model is accurate even if: (i) such features are significantly cropped due to the limited swath width of the SAR, (ii) the features are partly covered by sea ice and (iii) land is covering significant parts of the images. By evaluating the model performance on multiple input image resolutions (pixel sizes of 500m, 1km and 2km), it is found that higher resolution yield the best performance. This emphasises the potential of using high resolution sensors like SAR for detecting polar lows, as compared to conventionally used sensors such as scatterometers.

Index Terms—Polar lows; Mesocyclones; Deep learning; SAR

I. INTRODUCTION

Polar lows belong to the class of mesoscale maritime cyclones (from now on referred to as mesocyclones) that form at high latitudes, typically due to cold air outbreaks from sea ice or snow covered regions [1]. They are characterised by rapid development, small scale, strong winds and heavy snowfall. This makes them both difficult to predict and extremely hazardous for maritime activities such as fishing, shipping, petroleum extraction, and offshore wind power production. When making landfall, polar lows are prone to disrupt land and air traffic, destroy infrastructure, and trigger high snow avalanche activity in mountainous regions.

Due to their unpredictable and destructive nature, reliable and precise methods for early detection and tracking of polar lows are desirable. Meteorologists and scientists largely rely on direct observations in terms of satellite imagery or numerical weather prediction (NWP) models constrained by observations for detecting polar lows [2]–[9]. In maritime and polar regions, observations almost exclusively originate from

satellites. Conventionally, data from scatterometers, radiometers and optical sensors are assimilated into the NWP models [10], [11]. However, these sensors either rely on sunlight or have a coarse spatial resolution (typically a few to tens of kilometres). Considering that polar lows often occur during the polar night and are small scaled, featuring wind streaks, sharp atmospheric fronts and precipitation cells, observations at higher resolution regardless of light conditions could be beneficial.

Synthetic aperture radars (SARs) are independent of solar illumination and provide imagery at very high spatial resolution (typically a few to tens of metres). Researchers have already indicated that SAR data adds value to polar low monitoring [12]–[14]. Assimilation of SAR data into NWP models is however challenging, since the exact relationships between radar measurement and geophysical parameters are not trivial, especially at high wind speeds [14]–[16]. An alternative approach to make use of SAR data is to rely on data driven techniques, such as deep learning.

Deep learning has successfully been applied to several remote sensing applications and achieved state of the art results [17]–[20]. Cyclone type phenomena specifically, has been considered in assimilated data [21]–[23], passive microwave data [24], thermal infra-red (IR) data [25]–[29] and scatterometer data [30]. With the exception of [31], deep learning has however been largely overlooked for detecting mesocyclones in SAR data.

This paper investigates the possibility of using deep learning for detecting mesocyclones in general, and polar lows in particular, in SAR images. We aim to answer two main questions: (i) can a deep learning model recognise polar lows in SAR images, and (ii) what significance does the image resolution have on the performance?

To answer these questions, we first show that a training dataset can be constructed from the Sentinel-1 data archive, which is large enough for a deep neural network to be trained. In order to make the dataset large enough, we relax the definition of a polar low to the broader class of mesocyclones. The constructed dataset contains image samples divided in two classes, representing the presence and the absence of mesocyclones, respectively. In the following, we explain in detail how the dataset is built. To our knowledge, it is the first of its kind to be publicly released.

Then, we show how a deep neural network can be trained on the dataset to perform binary classification with very good performance. The deep learning model and the training procedure is carefully motivated by considering the training dataset size, input image size, and class imbalance. The performance of

*jgra@norceresearch.no

J. Grahn is with NORCE, The Norwegian Research Centre AS

F. M. Bianchi is with the Dept. of Mathematics and Statistics, UiT the Arctic University of Norway and with NORCE, The Norwegian Research Centre AS

the model is evaluated for multiple input image resolutions and interpretability techniques are applied on the model to evaluate what image features are most relevant for the classification.

II. SAR DATASET

This section describes the construction of the dataset for classifying mesocyclones, observed by the Sentinel-1 satellites. The dataset is publicly available (<https://doi.org/10.18710/FV5T9U>) and consists of 2004 images divided in two classes: the positive class (318 images with mesocyclones) and the negative class (1686 images without mesocyclones).

A. Positive class: Mesocyclone present

To build the positive class, polar lows monitored by the Sentinel-1 satellites were required. Historic catalogs of polar lows exist [32], [33], based on manual analysis of NWP model data as well as satellite data (thermal infrared, passive microwave and scatterometer data). However, these catalogs are regional and, more importantly, do not cover the time period when the Sentinel-1 satellites were operational. On the other hand, studies like [2], [3], [34] proposed objective criteria based on meteorological parameters that produce results similar to the manually annotated catalogs. Such objective criteria can be applied on reanalysis data, enabling identification of candidate low pressures that were coincident with the Sentinel-1 satellites.

Although a variety of objective criteria have been proposed, they are typically associated to either: (i) the low pressure intensity, (ii) the presence of a cold air outbreak, or (iii) the location of the low pressure in relation to the polar front. In [2], a combination of such criteria were imposed on the ECMWF reanalysis Interim (ERA-I) dataset and the most effective criteria for detecting polar lows were identified using the manual catalog by [32] as reference. However, events meeting all criteria are infrequent, since polar lows are rare. For reference, in [32], only 12 events per year were recorded over the Nordic seas on average, from year 2000 to 2009. Moreover, considering the limited spatio-temporal coverage of the Sentinel-1 satellites, not all events are imaged, making the number of image candidates even lower. Therefore, to include as many observed events as possible in our dataset, the cold air outbreak and locality type criteria were neglected. By considering only an intensity criteria, mesocyclones that are not necessarily driven by baroclinic instabilities or located in the polar air masses were included. Assuming that such mesocyclones share substantial similarities to polar lows, they can still provide valuable information to train a deep learning model, which motivates their inclusion in the dataset.

The intensity criteria was imposed on the ERA5 dataset. Specifically, it was formulated in terms of the depression in the sea level pressure (SLP) relative to the local mean. This type of criteria was considered by [2], where different SLP depression thresholds were tested. In our study, the threshold was set at 230 Pa and the local mean was computed within a 9×9 grid cell neighbourhood (corresponding to 270×270 km at the equator). The spatio-temporal distributions of resulting candidates and the subsequent matched SAR observations are

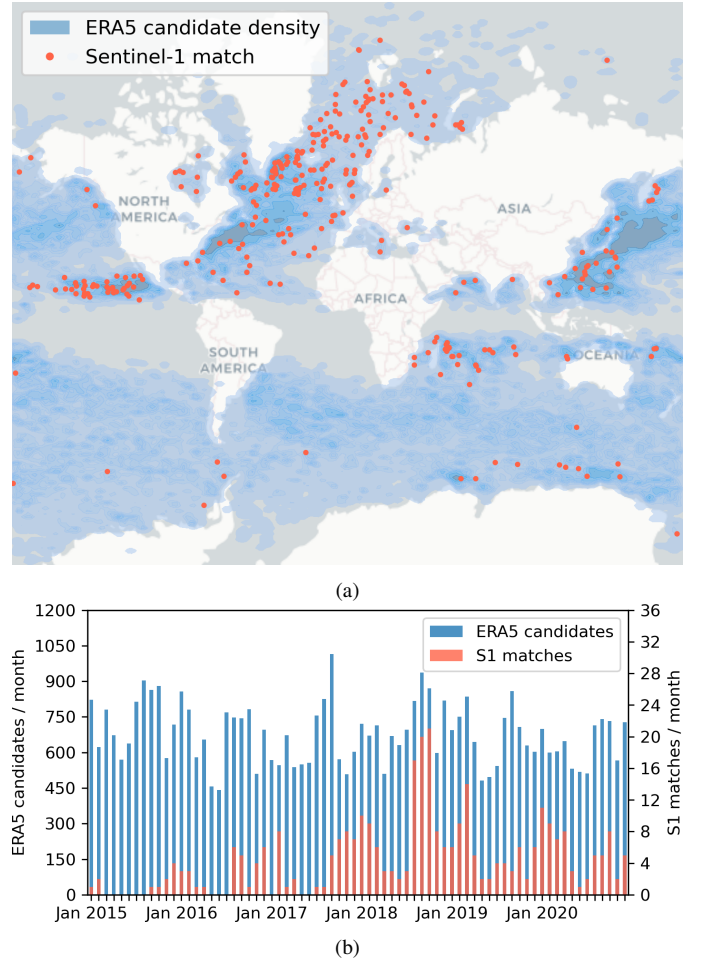


Fig. 1: The spatial distribution of ECMWF reanalysis version 5 (ERA5) candidates and Sentinel-1 matches in (a). The corresponding temporal distribution, as counts per month, in (b). Background map: © OpenStreetMap contributors/CARTO.

shown in figure 1. The highest concentrations of candidates were found in the subtropical regions of the North Pacific and North Atlantic. However, due to higher satellite revisit frequencies at higher latitudes, most SAR observations were found in the extra tropical and polar parts of the North Atlantic.

The dataset construction process is illustrated in figure 2 and each step is described in detail below.

1) *ERA5 filtering*: The ERA5 dataset [11] consists of hourly reconstructions of a large number of meteorological variables, spanning from 1950 to present. The data can be accessed on a geodetic grid, with a grid spacing of 0.25 degrees horizontally. Considering the global grid, candidate low pressures were identified by: (i) low-pass filtering the SLP using a 9×9 sliding average filter, (ii) selecting candidate grid cells where the SLP was 230 Pa lower than the low-pass SLP, (iii) grouping adjacent candidate grid cells and (iv) keeping groups with an equivalent radius smaller than 200 km (i.e. with an area less than $200^2 \pi$ km², thus excluding very big weather systems). Each such group was vectorised and constituted a candidate area of interest (AOI). This filtering process was done from 1 January 2015 to 31 December 2020 with a time

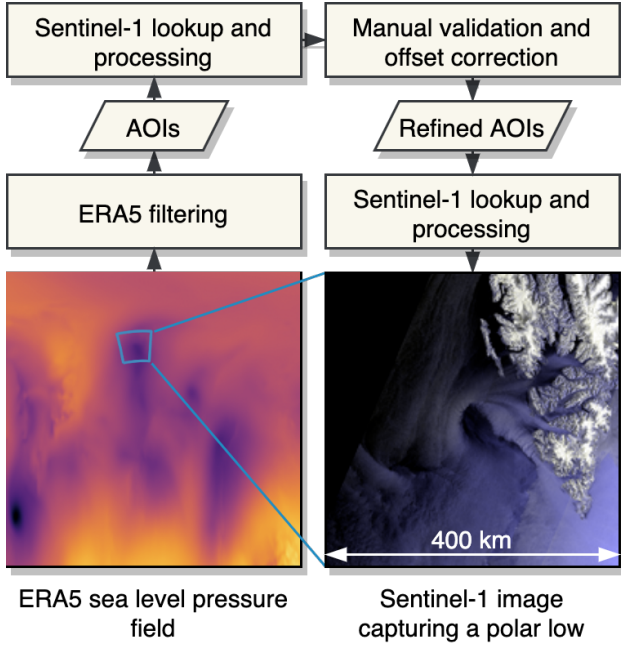


Fig. 2: The dataset construction process illustrated as a flow chart. The starting point is the ERA5 SLP field that is filtered according to our working definition of a polar low to form candidate AOIs. These are used to search for Sentinel-1 data, which is processed (see section II-A2) onto a UTM-grid. The processed images are validated manually and offset corrected by refining the AOIs and reprocessing the images onto a UTM-grid in which the image features are centered. The polarisation channels for this image sample is HH/HV.

step of 3 hours.

2) *Sentinel-1 lookup and processing*: For each candidate AOI at time t , we queried the Copernicus Open Access Hub for Sentinel-1 ground range detected (GRD) products in the time interval $t \pm 1.5$ hours. Resulting products were downloaded from Alaska satellite facility (ASF) and processed¹ by: (i) calibrating the data to sigma-nought, (ii) removing thermal noise, (iii) merging time-adjacent products to a common grid in SAR geometry, (iv) multi-looking to 500 m resolution in range and azimuth², (v) geocoding to a 400×400 km grid (centered at the AOI) in a universal transverse mercator (UTM) coordinate system with a 500 m grid spacing and (vi) generating red-green-blue (RGB) colour composites.

The RGB colour composites were generated by first re-scaling the radar cross-section (assumed in decibel scale) to a value $x \in [0, 1]$. Specifically, the 2nd and 98th percentiles of each separate image and polarisation channel were re-scaled to 0 and 1, respectively³.

¹All SAR data processing was done using Generic DAta Raster (GDAR), a python based library for processing raster data in radar geometries, developed by Norwegian Research Center (NORCE)

²In terms of number of looks, EW mode products are in total multi-looked by 60×20 looks in range and azimuth, while IW mode products are multi-looked by 250×50 looks in range and azimuth. Speckle noise is thus significantly suppressed in the processed images.

³The 2nd percentile was clipped to the range -25 to -15 dB and the 98th percentile was clipped to the range -10 to 0 dB. The clipping values were chosen to harmonise the scaling across image samples. Pixels without data were excluded when computing the percentiles and replaced by zeros.

For data with dual polarisation channels, the re-scaled values were used to make RGB colour composites as:

$$R = G = \frac{x_{||} + x_{\times}}{2}, \quad B = x_{||}$$

where $x_{||}$ and x_{\times} corresponds to the co- and cross-polarised⁴ channels, respectively. For single polarisation data, containing only the co-polarised channel, the colour channels were defined as: $R = G = B = x_{||}$. Both dual and single polarisation data were thus considered jointly in the training data set⁵, however, the dual polarisation data constituted the great majority of the samples (see figure 5).

3) *Manual validation and offset correction*: Each RGB colour composite was manually validated. Specifically, in each positive image, we asserted the presence of distinctive features (typically an eye or a comma shaped pattern). In general, these features were not centered in the processed images, since the image grid was centered at the candidate AOI originating from ERA5. Therefore, offsets were corrected for by manually centering the AOIs on the eye or comma shaped pattern. The samples were then reprocessed with the refined AOIs.

B. Negative class: Mesocyclone absent

To obtain samples of the negative class, representing the absence of a cyclone, we considered repeat-pass SAR acquisitions successive to those of the positive image samples (i.e. images acquired at the same relative orbit).

The motivation of our choice was twofold: (i) Sentinel-1 repeat-passes are separated by at least 6 days, which is enough time for the sea state (and thus the image features) to decorrelate, and (ii) the imaging geometry of repeat-pass acquisitions is nearly identical, such that static/background features appear similar. The second point is important in order to factor out land features from the dataset. Indeed, if the same land features appear in both the positive and negative class, it is expected that a machine learning model will be able to ignore them in the classification task.

As an example, a repeat-pass image set consisting of one positive and eight negative samples is shown in figure 3. To the left, the processed RGB composites are shown. The south tip of Svalbard can be seen statically in all images, while ocean features appear dynamically. The positive sample in the centre, contains a distinct vortex structure. To the right, a map with the footprints of the individual Sentinel-1 products involved is shown, together with the footprint of the image grid. Typically, due to the limited swath of the SAR, the products do not cover the whole image grid across track, leading to missing data in the RGBs in the cross-track direction. Occasionally, some products are not captured, leading to missing data in the along-track direction as well.

The distribution of SLP depression, defined as the difference in SLP (extracted from ERA5) between the image wide average and the average of the centre 100×100 pixels, is shown in figure 4. The positive samples have a strong depression,

⁴The co-polarised channel can be either HH or VV, and the corresponding cross-polarised channel can be either HV or VH.

⁵A dedicated experiment using only the co- or cross-polarised channel separately, can be found in the supplementary material.

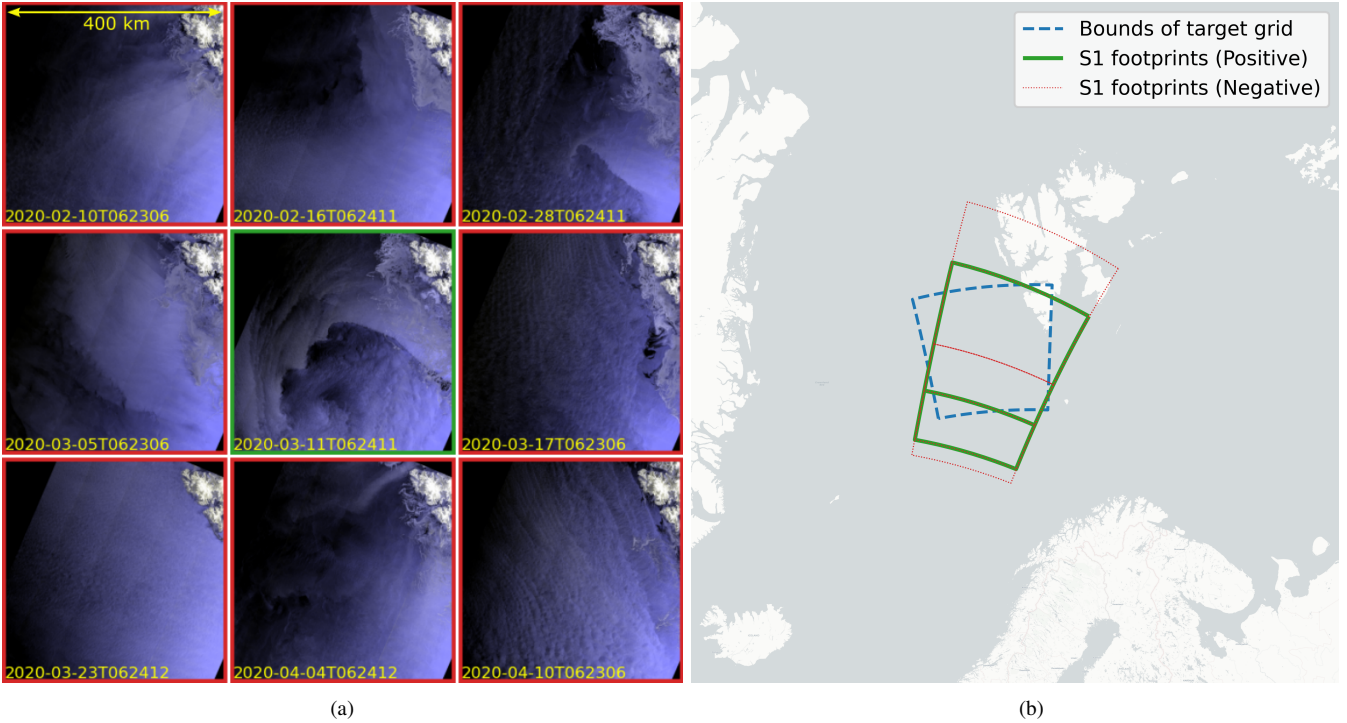


Fig. 3: In (a), a repeat-pass set is shown, consisting of one positive (central image) and eight negative samples. All samples within a set are processed onto the exact same grid, centered at the low pressure in the positive sample. The polarisation channels for this set are HH/HV. In (b), the location of the individual products and the grid is displayed. Background map: © OpenStreetMap contributors/CARTO.

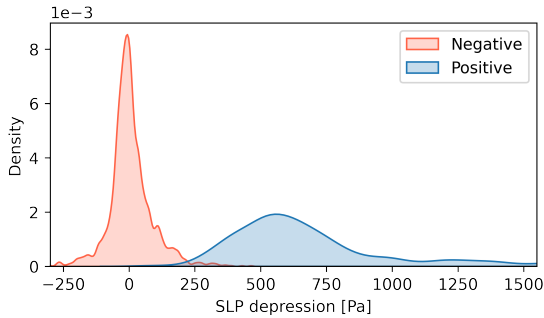


Fig. 4: The distribution of SLP depression for the two classes. The depression is measured as the SLP averaged over the whole image, minus the SLP averaged over the centre 100×100 pixels.

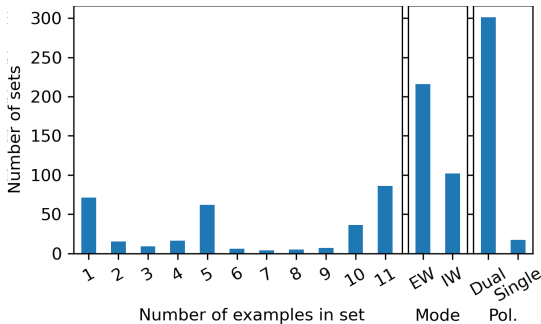


Fig. 5: Overview of the repeat-pass sets, in terms of distributions in size (samples per set), imaging mode and polarisation mode. For each positive sample, we extracted a maximum of 10 negative samples.

while the negative samples exhibit a symmetric distribution centred around 0 Pa. This indicates that the negative class, indeed, represents a sea state not biased towards a centered low

pressure. It should however be emphasised that no particular features are excluded from the negative class.

In total, 1686 negative samples were generated from the 318 positive ones, resulting in a total of 2004 samples in the dataset. The number of negatives per positive varied depending on the existence of repeat-pass acquisitions, as shown in figure 5. Furthermore, most samples were acquired in the extended wide-swath (EW) mode with two polarisations.

III. DEEP LEARNING

Three immediate challenges can be identified when choosing and training an appropriate deep learning model to perform classification on the dataset: (i) the input image size is relatively large, (ii) the training dataset is relatively small, and (iii) the classes are imbalanced. In the following, we discuss how these were dealt with.

A. Deep learning architecture

One of the major benefits of using SAR data, compared to e.g. scatterometer or passive microwave data, is the high image resolution. Although the images in the training dataset were already heavily downsampled from the original resolution of 10-40 m to 500 m, resulting in an image size of 800×800 pixels, they are relatively large in the context of many popular deep learning models. These are often designed for images of size 256×256 pixels or lower. To preserve details that are specific for the SAR data, such as wind streaks, rain cells or sharp atmospheric fronts, and to enable us to study the added value of high input image resolution, we wish to avoid further downsampling and rather let the model handle the

high input image resolution. Convolutional neural networks (CNNs) used for image classification usually consist of a stack of convolutional layers followed by pooling. Each such processing block sequentially increases the feature dimensionality through colvolutions, while reducing spatial resolution through pooling. As such, relevant spatial information will gradually become embedded in the feature space. If the input image is large, the model must either apply an aggressive downsampling in each processing block, or include many blocks and, thus, become very deep. The former can be obtained by large stride in the convolutional and pooling layers, or by using Atrous convolutions [35]. These techniques do, however, come at a cost of discarding spatial information, which we wish to avoid. This leaves us with the option of using a deep architecture, which gradually distill the spatial information and embeds it into the feature space.

Training a very deep network poses two fundamental challenges. Firstly, the gradients of the loss used to update the parameters may gradually vanish as they are backpropagated through the network. Secondly, an architecture with many layers contain many trainable parameters. This makes the model prone to overfitting, unless the training set is exceptionally large, which was not the case in our study.

A solution to address the first problem is to use residual connections, popularized by architectures such as ResNet [36], which facilitate the flow of the gradients during the backpropagation.

Considering the second problem, a ResNet is unfortunately characterized by many trainable parameters. There are, however, more recent deep architectures which include residual connections but have fewer parameters. For example, MobileNet [37] and Xception [38] implement separable 2D convolutions (Sep2DConv), which allows to greatly reduce the number of trainable parameters⁶.

Therefore, we opted for a customized Xception architecture⁷, whose details are depicted in figure 6. The entry block consists of a convolutional layer, followed by a batch normalization layer [39] and a ReLU activation function. There are L residual blocks, each one including Sep2DConv layers, batch normalization, ReLU activations and a max-pooling layer. The max-pooling output is combined with the input of the residual block through a skip connection. The convolutional layer in the middle of the skip connection has no activation function and simply applies a kernel of size 1 with stride 2, to match the shape of the input with the one of the output. A global pooling layer reduces the feature map generated by the last residual layer to a single vector, which is processed by the final classifier consisting of a dropout layer [40], a dense layer, and a softmax activation.

⁶The basic idea behind a Sep2DConv is to replace a matrix of parameters $\mathbf{W} \in \mathbb{R}^{N \times M}$ with an outer product of two unidimensional vectors, $\tilde{\mathbf{W}} = \mathbf{u}^T \otimes \mathbf{v}$, where $\mathbf{u} \in \mathbb{R}^M$ and $\mathbf{v} \in \mathbb{R}^N$, reducing the number of parameters from $M \cdot N$ to $M + N$. This is, actually, a simplification. In practice, a Sep2DConv splits the traditional convolution with a kernel of size $H \times W \times F_{in} \times F_{out}$ with a depth-wise convolution with F_{in} kernels of size $H \times W \times 1$, followed by a point-wise convolution with a kernel of size $1 \times 1 \times F_{out}$.

⁷A comparison with other popular deep learning architectures is presented in the supplementary material.

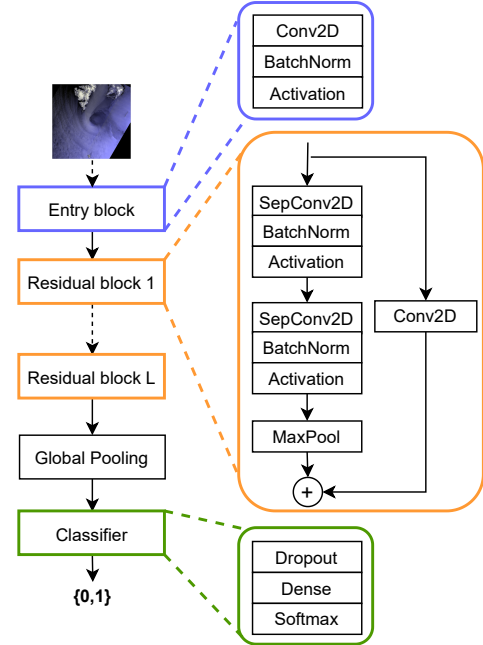


Fig. 6: Architecture details.

B. Data augmentation

Augmenting the training data by applying random transformation is a common technique used to prevent overfitting. By exposing the deep learning model to perturbations of original inputs, it is possible to improve the robustness of the model. In addition, data augmentation allows to get rid of some bias in the dataset and increase the generalization performance on new unseen data.

Our dataset has been designed by keeping image augmentation in mind. Each low pressure is centered in the image and has a wide area around that can be partially cropped. Each time a batch of images is fetched to our deep learning model, the following random transformations are applied on the fly; (i) horizontal and vertical translation (between 0 and 10% of the image size), (ii) horizontal and vertical flip, (iii) rotation (0 to 40 degrees), (iv) zoom (-10% to 10% of the original scale) and finally (v) cropping to the centre 512×512 pixels. If after the transformation some points fall outside the boundaries of the original input image, these are filled with zeros. Notably, after data augmentation the low pressures are no longer centered in each image. Figure 7 shows an example of augmented images randomly generated during training.

C. Class imbalance

While the number of positive samples were restricted by the number of matches found on the Sentinel-1 archive, multiple negative samples could be generated for each positive sample. This lead to a natural skewness in the distribution of classes in the dataset: 84% of the samples belong to the negative class and 16% to the positive. We tested and compared three different approaches to train the deep learning model in the presence of class imbalance.

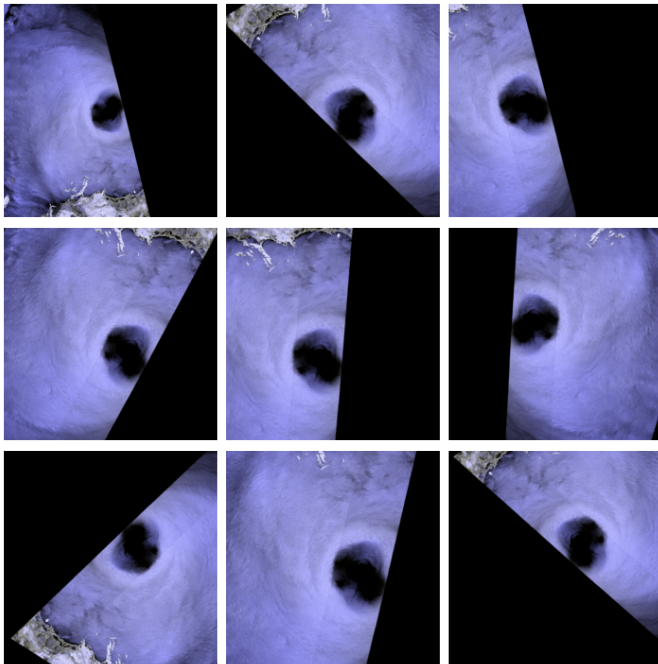


Fig. 7: Examples of random image augmentation. The top-left image is the original. The polarisation channels for the image are VV/VH.

1) *Class-weighting*: We re-weighted the loss function according to class frequencies. Denoting the number of samples in the negative and positive classes respectively by n_0 and n_1 , the loss for samples in the corresponding classes were weighted by $w_0 = \frac{1}{n_0} \frac{n_0+n_1}{2}$ and $w_1 = \frac{1}{n_1} \frac{n_0+n_1}{2}$. This means that each error on the positive class affects the optimization of the model weights to a greater extent.

2) *Oversampling*: We ensured that in each batch there were always the same amount of samples of both classes. Specifically, this was obtained by designating half of the batch to the negative class, and half to the positive. In each epoch⁸, samples from the negative class are seen only once, while samples from the positive class are repeated. We made sure to observe all samples of the negative class at least once during an epoch and to not sample any image twice within a batch.

3) *Rejection sampling*: This strategy drops samples until a balanced distribution across the two classes is obtained. Contrarily to oversampling, each sample is seen at most once in each epoch, which makes the overall training faster.

By empirical comparison (see the supplementary material), we found that oversampling yields the best performance and, therefore, was the strategy adopted in our experiments.

D. Hyperparameter tuning

To find the optimal configuration of the deep learning model, we searched several hyperparameters and selected those giving the best performance on a validation set. As validation set, we used 10% of the training set. The hyperparameter space and the optimal values found after the optimization procedure are reported in table I. To reduce the

⁸One epoch is when the whole training dataset has been passed forward and backward through the network once.

hyperparameters space, we only search the number of filters of the first residual blocks and then we double the number in the following blocks.

Hyperparam.	Search space	Optimal
Activation	{ReLU, SeLU}	ReLU
Conv2D filters	{8, 16}	8
Kernel size	{3, 5}	3
SepConv2D filters	{8, 16, 24, 32}	8
Num. residual blocks	{2, 8}	7
Global pooling	{avg, flat, max}	avg
Num. dense layers	{1, 3}	1
Units in the dense layer	{8, 16, 24, 32}	8
Dropout rate	{0.1, 0.6}	0.5
Use batch normalization	{True, False}	True
Learning rate	{1e-2, 1e-3, 1e-4}	1e-3

TABLE I: Hyperparameters space and optimal values found. “Conv2D filters” and “kernel size” refer to the entry block. “Sep-Conv2D filters” refers to the 1st residual block, since the number of filters is double each time in the following blocks.

Since the dataset contains large images and we consider deep models with many parameters, evaluating each hyperparameter configuration is computationally expensive. Therefore, rather than performing an exhaustive search with grid search or evaluating a large number of configurations with a random search, we opted for a more efficient approach. In particular, we used Bayesian hyperparameter optimization [41].

We used a batch size of 16 and the Adam optimizer [42]. During the hyperparameter tuning we trained the model for 50 epochs. After finding the optimal configuration, we trained the final model for 200 epochs.

E. Model interpretability

Due to the presence of many non-linear transformations, it is difficult to interpret the decision process of a neural network and considerable research effort has been devoted to improve our understandings. Gradient based approaches try to find which inputs have the most influence on the model scoring function for a given class [43]–[45]. This is usually done by taking the gradient of the class activation score with respect to each input features [46]. A drawback of gradient based methods is that they give zero contribution to inputs that saturate the ReLU or MaxPool. To capture such shortcomings, a formal notion of explainability was introduced in [47] with the axiom of conservation of total relevance, which states that the sum of relevance of all pixels must match the class score of the model. Specifically, the authors propose to distribute the total relevance of the class score to the input features with Layer-wise Relevance Propagation (LRP). While LRP follows the conservation axiom, it does not specify how to distribute the relevance among the input features. To address this problem DeepLiFT [48] enforces an additional axiom on how to propagate the relevance by following the chain rule.

In this work, we adopt two recent interpretability techniques, that address some of the shortcomings discussed above and are able to provide valuable insights into the decision problem of our model.

1) *Integrated Gradients*: Integrated gradients (IG) [49] has become a popular interpretability technique since it can be applied to any neural network model, is easy to implement,

and theoretically grounded. IG aims to satisfy two additional axioms that are not jointly ensured by other existing attribution schemes; (i) if the input and an uninformative baseline differ in exactly one feature, such a feature should be given non-zero attribution, (ii) when two models are functionally equivalent, they must have identical attributions to input features.

Denoting the model scoring function F , the attributions given by IG are

$$\text{IG}(x) := (x - x') \cdot \int_{\alpha=0}^1 \frac{\partial F(x' + \alpha \cdot (x - x'))}{\partial x} d\alpha, \quad (1)$$

where x is a sample in the dataset, x' is the uninformative baseline, and α is an interpolation constant used to perturb the input features.

In our study, we let x' be a black image (all zeros) as the uninformative baseline. As empirically confirmed in our experiments, such a baseline is classified with high confidence to be negative. Let \mathcal{X} be the set of interpolated images from x' to x . The computation of the integral in (1) is approximated with the sum of the partial derivatives of the images in \mathcal{X} . Figure 8 depicts a small interpolation set \mathcal{X} from the mean-baseline to a positive sample and shows how the classification score changes. By summing the gradients $\frac{\partial F(\mathcal{X})}{\partial x_i}$ one quantifies

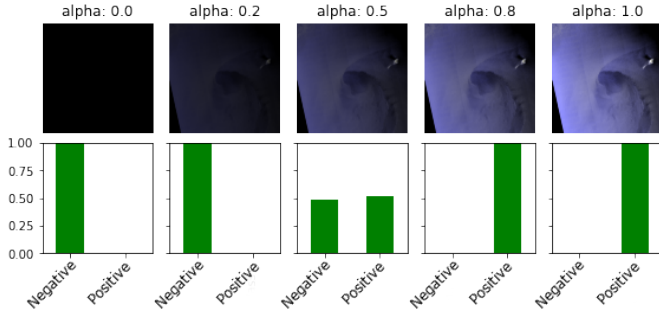


Fig. 8: Top row: linear interpolation from the zero-baseline (left) to an actual sample of positive class (right). Bottom row: classification probabilities assigned by the model at each step of the interpolation.

the relationship between the changes in the input features and the changes in the predictions of the model F .

2) *Gradient-weighted Class Activation Mapping*: While IG can be used on any neural network model, gradient-weighted class activation mapping (Grad-CAM) is specific for CNNs. It uses the gradients of a given target class flowing into the final convolutional layer to produce a coarse localization map, which highlights the important regions in the image for predicting the class [50]. We summarize at a high level the main steps of the algorithmic procedure and we refer the interested reader to [50] for further details: (i) take a trained model and cut it at the k -th layer, which is the layer for which we want to create a Grad-CAM heat-map (usually, the activation after the last convolutional layer), (ii) feed an input image to the complete model and collect the total loss and the output of layer k , (iii) compute the gradients of the output of layer k with respect to the loss, (iv) take parts of the gradient which contribute to the prediction and use to build a heatmap, and (v) resize the heatmap so that it can be overlaid to the original image.

IV. RESULTS

The dataset presented in section II was used to train the model described in section III. Specifically, the dataset was partitioned such that 79% of the samples were used for training and validation and 21% for testing. The partitioning was done by randomly assigning complete repeat-pass sets to either the test or the training set. In such a way, positive and negative samples with the same land features cannot appear both in the training and test set. This, (i) encouraged the model to factor land features out as irrelevant to the classification task, (ii) allowed us to evaluate the generalization capability of the model by testing on new locations, unseen during training.

The arguably most attractive property of SAR data, as compared to e.g. scatterometer data, is the high spatial resolution. In order to evaluate the added value of higher spatial resolution, the model accuracy was examined for three different input image resolutions⁹; 500 m, 1000 m and 2000 m (the latter two obtained by bi-linear down sampling of the first). Hyperparameter tuning was performed independently for each resolution (see the supplementary material for details), and the classification performance on the test set is shown in table II. The table displays the mean and standard deviation of true negatives (TNs), false negatives (FNs), false positives (FPs), true positives (TPs) and F1 score obtained from 10 independent runs. It is clear that higher image resolution significantly improves the classification results¹⁰. In fact, for the highest input resolution, the model is misclassifying on average less than 8 samples (as FN or FP) out of the 435 samples in the test set, with a mean F1 score of 0.94 (in the supplementary materials, results of the performance for the highest input image resolution using the co- or cross-polarised channels separately are also presented).

Pixel size	TN	FN	FP	TP	F1 score
2km	346.6±2.1	6.4±2.1	9.8±1.9	54.2±1.9	0.87±0.01
1km	364.4±2.1	6.6±2.1	7.4±2.1	56.6±2.1	0.89±0.02
500m	367.8±2.7	3.2±2.7	4.6±1.7	59.4±1.7	0.94±0.01

TABLE II: Classification performance on the test set when using different input resolutions. It is evident that higher input image resolution significantly improves the performance.

TN	FN	FP	TP	F1 score
366	2	5	62	0.95

TABLE III: Classification performance obtained on the specific run where we apply the interpretability techniques.

A model trained on the 500 m resolution images was further examined using the IG and Grad-CAM techniques presented in section III-E. The performance of this specific model is shown in table III and the images it classifies as TPs, FPs, and FNs are discussed in the following. The deep learning model used

⁹The highest resolution here (500 m) is still considerably lower than the original resolution of the SAR images. However, as discussed in section III-A, the input image size is limited by the depth of the network architecture in relation to the size of the training dataset. Therefore, we did not considered even higher input image resolutions, even if the original data allowed for it.

¹⁰A detailed comparison based on Grad-CAM between the model trained on 2000 m and 500 m resolution is presented in the supplementary materials.

in our experiments and the code to apply the interpretability techniques is available online¹¹.

A. True positives

Of the 62 TP samples (i.e. low pressures correctly classified as low pressures), 4 samples are displayed in figure 9. The first column shows the input RGB colour composites, the second column shows the IG in green and the third column shows the Grad-CAM as a heat map. 3 out of these 4 samples are located in polar regions, while the sample on the second row is an extra-tropical cyclone observed off the coast of Japan. The IG and Grad-CAM overlays indicate that the model is focusing on the cyclonic eye features. The IG overlay has a slightly higher emphasis on the wind fronts as compared to the Grad-CAM. Both the IG and Grad-CAM indicate that the model is effectively disregarding land features as well as the sea ice features appearing in the top row. Notably, in the top row, a large part of the cyclonic eye feature is also cropped due to the limited swath width of the SAR. This is the case in multiple samples classified as TP, indicating a certain robustness to image features being cropped or obscured by e.g. sea ice.

B. False positives

The model classified 5 samples as FP (i.e. absence of low pressures incorrectly classified as low pressures), of which 4 are shown in figure 10. The top two samples are presumably difficult to classify correctly (or the ground truth label could potentially be wrong), as they actually contain some pronounced wind fronts. Considering the IG and Grad-CAM, indeed the model is focusing on these wind features. The sample on the third row also contains a pronounced wind front that the model is focusing on, but the front is not curved. The classification score is however only 0.57 for this sample. In the fourth sample, no wind front is visible, but the IG and Grad-CAM reveal that the model focuses on a wind wake (formed behind the Izu peninsula, Japan, located in the image centre), which may be misinterpreted as a cyclonic eye.

Finally, we notice that IG and Grad-CAM highlight different areas in the second and third image. Explainability techniques for deep learning are tools meant for diagnostic, which require a certain degree of subjective interpretation. Each technique is based on specific heuristics, which put a bias on what features are considered relevant. Indeed, even for samples classified with high confidence two explainability techniques might focus on different input features [51]. The discrepancy is often exacerbated in samples classified with lower confidence.

C. False negatives

Only two samples of the positive class were incorrectly classified as negatives, i.e. mistaken as absence of low pressure while being labeled as low pressures. These are shown in figure 11. Here, IG are not computed, since a black image cannot be used as a baseline for the negative class. Nevertheless, Grad-CAM can still be computed and is shown in

the second column. It can be noted that both images suffer from lacking data due to the limited swath width of the SAR acquisitions. Indeed, the Grad-CAM indicate that the model is not focusing on the darker center features as was the case for the TP samples in figure 9. It should however be emphasised that this happens for only 2 of the 368 negative samples in the test set.

V. CONCLUSIONS

In this study, we show that SAR images from the Sentinel-1 satellites provide an attractive data source for automatic and accurate detection of maritime mesocyclones, such as polar lows. Specifically, we show that sufficiently many image examples can be found to build a labeled dataset for a deep learning model to be trained. By further comparing our deep learning model when trained on different input image resolutions, we find that higher resolution yields significantly better performance. This highlights the added value of using SAR data, as compared to conventionally used sensors of lower resolution. In particular, at 500 meters resolution we get an F1 score of 0.94, as compared to 0.87 at 2 km resolution (comparable to modern scatterometers).

It should further be noted that the highest resolution tested in this study (500 meters) is primarily limited by the size of the training dataset and not the native resolution of the SAR sensor (10-40 meters). Thus, even higher input image resolutions could in principle be considered, potentially with even better performance. Larger input image sizes, however, ideally require deeper neural network architectures, with more trainable weights. This in turn require larger training datasets to avoid over-fitting. Even so, with an increasing amount of SAR data being available from new satellites, larger training datasets could be constructed in the future, enabling even better performance.

By design, the training dataset contains spiral-form low pressures in the positive class. By analyzing IG and Grad-CAM on the trained model, we verify that the spiral shaped atmospheric fronts and the low wind centres yield most of the class attribution. Moreover, we conclude that: (i) these characteristic wind features do not need to be fully covered in the images, but can be substantially cropped due to the limited swath width of the SAR, (ii) wind features can be partly covered by sea ice and still be identified by the model, and (iii) the model is able to ignore land features in the images. The last point can be verified thanks to the procedure used to obtain the negative samples, i.e. through repeat-pass acquisitions (see section II-B).

In summary, we conclude that the application of deep learning on SAR images for recognising maritime mesocyclones is promising. Further evaluation and comparison to detection based on data from other sensors or NWP models is encouraged as a future work direction.

VI. ACKNOWLEDGMENTS

This work was funded by the European space agency (ESA), under the open call project *Polar low detection based on Sentinel-1 data* (contract number 4000129961). We would

¹¹<https://github.com/FilippoMB/Recognition-of-polar-lows-in-Sentinel-1-SAR-images-with-deep-learning>

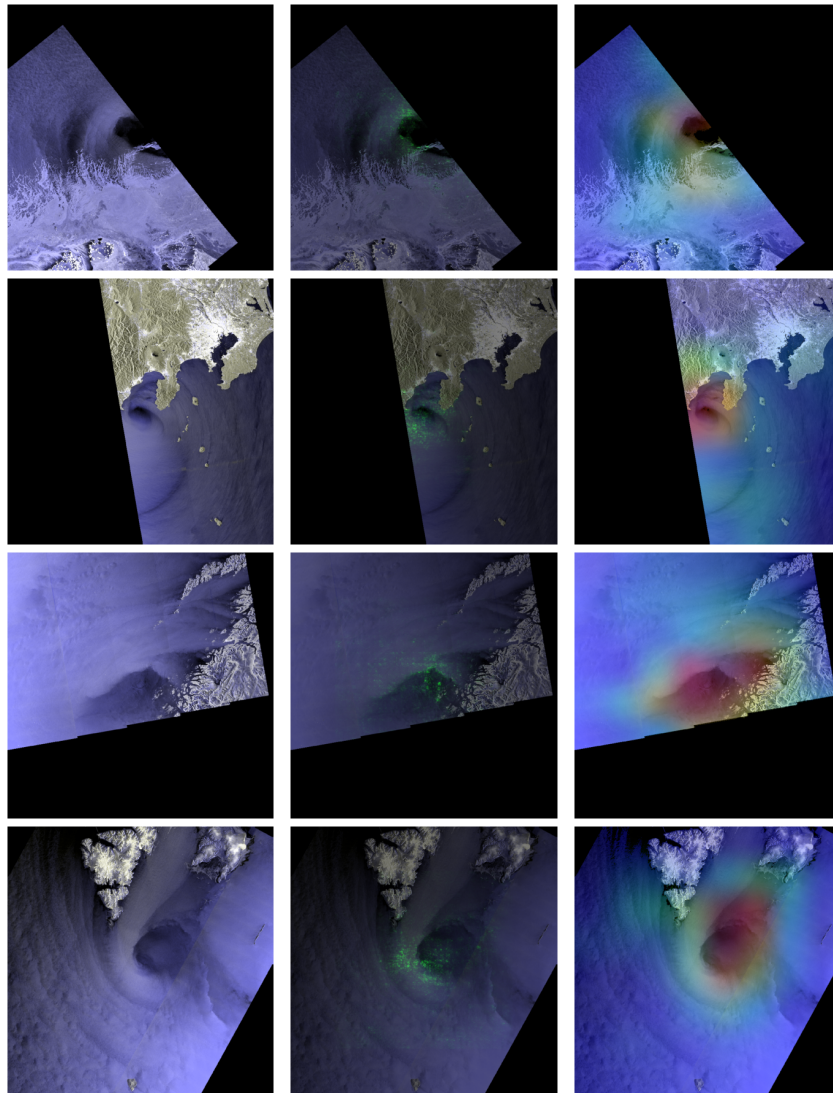


Fig. 9: True positives (4 of 62 samples shown) in first column and corresponding IG and Grad-CAM in the second and third column. The classification scores are 0.99, 0.82, 0.99, and 0.99, respectively. The polarized channels are HH/HV, HH/HV, VV/VH and HH/HV, respectively.

like to thank Patrick Stoll for his valuable feedback. We thank those involved in developing the GDAR software used to process SAR data, especially Heidi Hindberg, Yngvar Larsen, and Tom Grydeland. We also thank Temesgen Gebrie Yitayew and Hannah Vickers for their help in establishing this project. Finally, we would like to thank the reviewers for their insightful comments.

REFERENCES

- [1] E. A. Rasmusson, "Polar lows," in *A Half Century of Progress in Meteorology: A Tribute to Richard Reed*. Springer, 2003, pp. 61–78.
- [2] P. J. Stoll, R. G. Graversen, G. Noer, and K. Hodges, "An objective global climatology of polar lows based on reanalysis data," *Quarterly Journal of the Royal Meteorological Society*, vol. 144, no. 716, pp. 2099–2117, 2018.
- [3] W. Yanase, H. Niino, I. W. Shun-ichi, K. Hodges, M. Zahn, T. Spengler, and I. A. Gurvich, "Climatology of polar lows over the sea of japan using the jra-55 reanalysis," *Journal of Climate*, vol. 29, no. 2, pp. 419–437, 2016.
- [4] A.-M. Blechschmidt, "A 2-year climatology of polar low events over the nordic seas from satellite remote sensing," *Geophysical Research Letters*, vol. 35, no. 9, 2008.
- [5] K. Wilhelmsen, "Climatological study of gale-producing polar lows near norway," *Tellus A: Dynamic Meteorology and Oceanography*, vol. 37, no. 5, pp. 451–459, 1985.
- [6] G. Zappa, L. Shaffrey, and K. Hodges, "Can polar lows be objectively identified and tracked in the ecmwf operational analysis and the era-interim reanalysis?" *Monthly Weather Review*, vol. 142, no. 8, pp. 2596–2608, 2014.
- [7] C. Michel, A. Terpstra, and T. Spengler, "Polar mesoscale cyclone climatology for the nordic seas based on era-interim," *Journal of Climate*, vol. 31, no. 6, pp. 2511–2532, 2018.
- [8] L. Xia, M. Zahn, K. Hodges, F. Feser, and H. Storch, "A comparison of two identification and tracking methods for polar lows," *Tellus A: Dynamic Meteorology and Oceanography*, vol. 64, no. 1, p. 17196, 2012.
- [9] P. J. Stoll, T. Spengler, A. Terpstra, and R. G. Graversen, "Polar lows – moist-baroclinic cyclones developing in four different vertical wind shear environments," *Weather and Climate Dynamics*, vol. 2, no. 1, pp. 19–36, 2021.
- [10] M. Muller, M. Homleid, K.-I. Ivarsson, M. A. Ø. Køltzow, M. Lindskog, K. H. Midtbø, U. Andrae, T. Aspelien, L. Berggren, D. Bjørge, P. Dahlgren, J. Kristiansen, R. Randriamampianina, M. Ridal, and O. Vignes, "Arome-metcoop: A nordic convective-scale operational weather prediction model," *Weather and Forecasting*, vol. 32, no. 2, pp. 609 – 627, 2017.
- [11] H. Hersbach, B. Bell, P. Berrisford, S. Hirahara, A. Horányi, J. Muñoz-

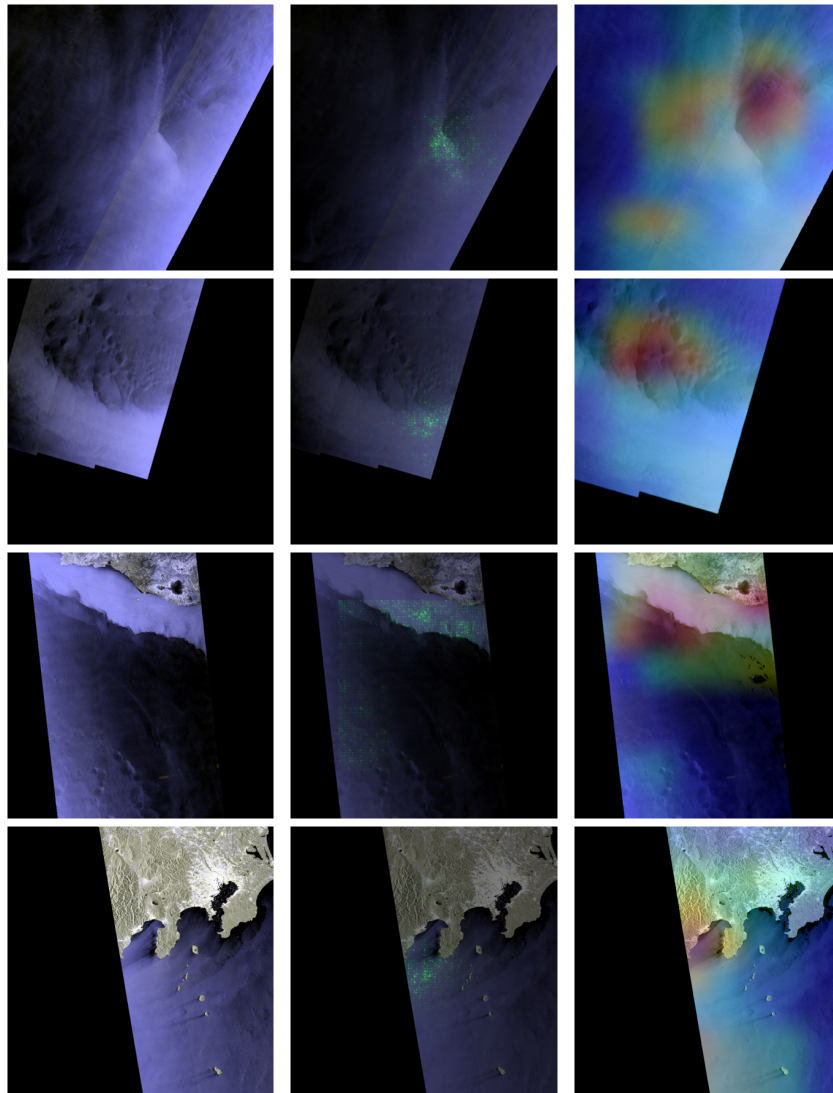


Fig. 10: False positives (4 of 5 samples shown), with classification scores 0.89, 0.92, 0.57 and 0.87. The polarised channels are HH/HV, HH/HV, VV/VH and VV/VH, respectively.

- Sabater, J. Nicolas, C. Peubey, R. Radu, D. Schepers *et al.*, “The era5 global reanalysis,” *Quarterly Journal of the Royal Meteorological Society*, vol. 146, no. 730, pp. 1999–2049, 2020.
- [12] G. W. K. Moore and P. W. Vachon, “A polar low over The Labrador Sea: Interactions with topography and an upper-level potential vorticity anomaly, and an observation by RADARSAT-1 SAR,” *Geophysical research letters*, vol. 29, no. 16, p. 1773, Aug. 2002.
- [13] B. R. Furevik, H. Schyberg, G. Noer, F. Tvetter, and J. Röhrs, “Asar and ascet in polar low situations,” *Journal of Atmospheric and Oceanic Technology*, vol. 32, no. 4, pp. 783–792, 2015.
- [14] M. Tollinger, R. Graverson, and H. Johnsen, “High-resolution polar low winds obtained from unsupervised sar wind retrieval,” *Remote Sensing*, vol. 13, no. 22, p. 4655, 2021.
- [15] B. Chapron, H. Johnsen, and R. Garello, “Wave and wind retrieval from sar images of the ocean,” in *Annales des telecommunications*, vol. 56, no. 11. Springer, 2001, pp. 682–699.
- [16] A. A. Mouche, F. Collard, B. Chapron, K.-F. Dagestad, G. Guitton, J. A. Johannessen, V. Kerbaol, and M. W. Hansen, “On the use of doppler shift for sea surface wind retrieval from sar,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 50, no. 7, pp. 2901–2909, 2012.
- [17] X. X. Zhu, D. Tuija, L. Mou, G.-S. Xia, L. Zhang, F. Xu, and F. Fraundorfer, “Deep learning in remote sensing: A comprehensive review and list of resources,” *IEEE Geoscience and Remote Sensing Magazine*, vol. 5, no. 4, pp. 8–36, 2017.
- [18] F. M. Bianchi, M. M. Espeseth, and N. Borch, “Large-scale detection and categorization of oil spills from sar images with deep learning,” *Remote Sensing*, vol. 12, no. 14, p. 2260, 2020.
- [19] F. M. Bianchi, J. Grahn, M. Eckerstorfer, E. Malnes, and H. Vickers, “Snow avalanche segmentation in sar images with fully convolutional neural networks,” *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 14, pp. 75–82, 2020.
- [20] L. T. Luppino, M. Kampffmeyer, F. M. Bianchi, G. Moser, S. B. Serpico, R. Jenssen, and S. N. Anfinsen, “Deep image translation with an affinity-based change prior for unsupervised multimodal change detection,” *IEEE Transactions on Geoscience and Remote Sensing*, 2021.
- [21] Y. Liu, E. Racah, J. Correa, A. Khosrowshahi, D. Lavers, K. Kunkel, M. Wehner, W. Collins *et al.*, “Application of deep convolutional neural networks for detecting extreme weather in climate datasets,” *arXiv preprint arXiv:1605.01156*, 2016.
- [22] D. Matsuoka, M. Nakano, D. Sugiyama, and S. Uchida, “Deep learning approach for detecting tropical cyclones and their precursors in the simulation by a cloud-resolving global nonhydrostatic atmospheric model,” *Progress in Earth and Planetary Science*, vol. 5, no. 1, pp. 1–16, 2018.
- [23] S. Giffard-Roisin, M. Yang, G. Charpiat, C. Kumler Bonfanti, B. Kégl, and C. Monteoloni, “Tropical cyclone track forecasting using fused deep learning from aligned reanalysis data,” *Frontiers in Big Data*, vol. 3, p. 1, 2020.
- [24] A. Wimmers, C. Velden, and J. H. Cossuth, “Using deep learning to estimate tropical cyclone intensity from satellite passive microwave imagery,” *Monthly Weather Review*, vol. 147, no. 6, pp. 2261–2282, 2019.

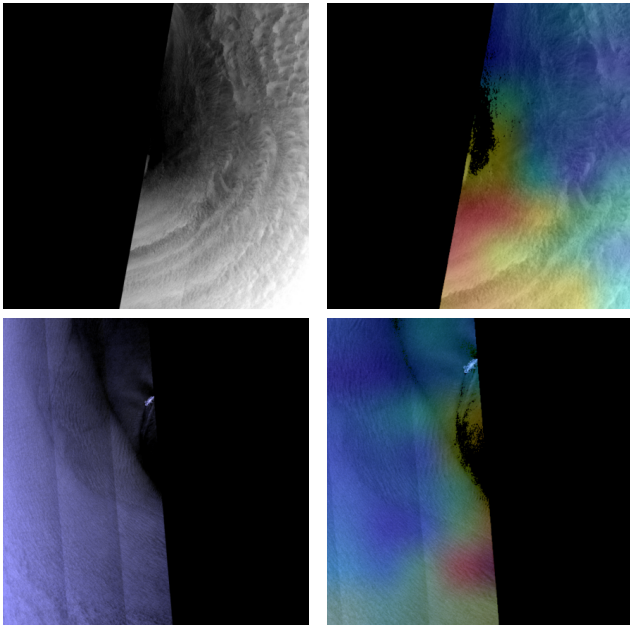


Fig. 11: False Negatives (all samples shown), with classification scores 0.96 and 0.79. The top image is acquired with the HH polarisation channel and the bottom image in the HH/HV polarised channels.

- 2019.
- [25] P. Golubkin, J. Smirnova, and L. Bobylev, "Satellite-derived spatio-temporal distribution and parameters of north atlantic polar lows for 2015–2017," *Atmosphere*, vol. 12, no. 2, 2021.
- [26] C. Kumler-Bonfanti, J. Stewart, D. Hall, and M. Govett, "Tropical and extratropical cyclone detection using deep learning," *Journal of Applied Meteorology and Climatology*, vol. 59, no. 12, pp. 1971–1985, 2020.
- [27] B.-F. Chen, B. Chen, H.-T. Lin, and R. L. Elsberry, "Estimating tropical cyclone intensity by satellite imagery utilizing convolutional neural networks," *Weather and Forecasting*, vol. 34, no. 2, pp. 447–465, 2019.
- [28] R. Pradhan, R. S. Aygun, M. Maskey, R. Ramachandran, and D. J. Cecil, "Tropical cyclone intensity estimation using a deep convolutional neural network," *IEEE Transactions on Image Processing*, vol. 27, no. 2, pp. 692–702, 2017.
- [29] M. Krinitskiy, P. Verezhenskaya, K. Grashchenkov, N. Tilinina, S. Gulev, and M. Lazzara, "Deep convolutional neural networks capabilities for binary classification of polar mesocyclones in satellite mosaics," *Atmosphere*, vol. 9, no. 11, p. 426, 2018.
- [30] M. Xie, Y. Li, and K. Cao, "Global cyclone and anticyclone detection model based on remotely sensed wind field and deep learning," *Remote Sensing*, vol. 12, no. 19, p. 3111, 2020.
- [31] A. R. Carmo, N. Longépé, A. Mouche, D. Amorosi, and N. Cremer, "Deep learning approach for tropical cyclones classification based on c-band sentinel-1 sar images," in *2021 IEEE International Geoscience and Remote Sensing Symposium IGARSS*. IEEE, 2021, pp. 3010–3013.
- [32] G. Noer, Ø. Saetra, T. Lien, and Y. Gusdal, "A climatological study of polar lows in the nordic seas," *Quarterly Journal of the Royal Meteorological Society*, vol. 137, no. 660, pp. 1762–1772, 2011.
- [33] J. E. Smirnova, P. A. Golubkin, L. P. Bobylev, E. V. Zabolotskikh, and B. Chapron, "Polar low climatology over the nordic and barents seas based on satellite passive microwave data," *Geophysical Research Letters*, vol. 42, no. 13, pp. 5603–5609, 2015.
- [34] T. J. Bracegirdle and S. L. Gray, "An objective climatology of the dynamical forcing of polar lows in the nordic seas," *International Journal of Climatology: A Journal of the Royal Meteorological Society*, vol. 28, no. 14, pp. 1903–1919, 2008.
- [35] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.
- [36] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Thirty-first AAAI conference on artificial intelligence*, 2017, pp. 1251–1258.
- [37] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [38] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.
- [39] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*. PMLR, 2015, pp. 448–456.
- [40] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [41] J. Snoek, H. Larochelle, and R. P. Adams, "Practical bayesian optimization of machine learning algorithms," *Advances in neural information processing systems*, vol. 25, 2012.
- [42] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *International Conference on Learning Representations*, 2014.
- [43] D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg, "Smoothgrad: removing noise by adding noise," *International Conference on Machine Learning*, 2017.
- [44] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *European conference on computer vision*. Springer, 2014, pp. 818–833.
- [45] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for simplicity: The all convolutional net," *International Conference on Learning Representations*, 2015.
- [46] K. Simonyan, A. Vedaldi, and A. Zisserman, "Deep inside convolutional networks: Visualising image classification models and saliency maps," *International Conference on Learning Representations*, 2014.
- [47] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek, "On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation," *PLoS one*, vol. 10, no. 7, p. e0130140, 2015.
- [48] A. Shrikumar, P. Greenside, and A. Kundaje, "Learning important features through propagating activation differences," in *International Conference on Machine Learning*. PMLR, 2017, pp. 3145–3153.
- [49] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in *International Conference on Machine Learning*. PMLR, 2017, pp. 3319–3328.
- [50] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.
- [51] W. Samek, G. Montavon, A. Vedaldi, L. K. Hansen, and K.-R. Müller, *Explainable AI: interpreting, explaining and visualizing deep learning*. Springer Nature, 2019, vol. 11700.
- [52] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015. [Online]. Available: <http://arxiv.org/abs/1409.1556>
- [53] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [54] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *CoRR*, vol. abs/1704.04861, 2017. [Online]. Available: <http://arxiv.org/abs/1704.04861>
- [55] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. [Online]. Available: <https://openreview.net/forum?id=YicbFdNTTy>
- [56] I. O. Tolstikhin, N. Houlsby, A. Kolesnikov, L. Beyer, X. Zhai, T. Unterthiner, J. Yung, A. Steiner, D. Keysers, J. Uszkoreit *et al.*, "Mlp-mixer: An all-mlp architecture for vision," *Advances in Neural Information Processing Systems*, vol. 34, pp. 24 261–24 272, 2021.
- [57] G. Zhang, X. Li, W. Perrie, P. A. Hwang, B. Zhang, and X. Yang, "A hurricane wind speed retrieval model for c-band radarsat-2 cross-polarization scansar images," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 8, pp. 4766–4774, 2017.

Supplementary material

I. COMPARISON WITH OTHER DEEP LEARNING ARCHITECTURES

Here, we report the results achieved with off-the-shelf deep learning architectures for image classification. In Tab. IV we report the results obtained with VGG16 [52], ResNet50 [53], Xception [38], MobileNet [54], ViT [55], and MLP Mixer [56]. The implementation of VGG16, ResNet50, Xception, and MobileNet is the one from *Keras applications*¹².

Architecture	TN	FN	FP	TP	F1 score
VGG16	371	64	0	0	0.0
ResNet50	366	5	11	53	0.87
Xception	365	4	6	60	0.92
MobileNet	369	11	2	53	0.89
ViT	350	9	21	55	0.79
MLPMixer	306	25	65	39	0.46

TABLE IV: Results obtained with popular architectures. The best performance in terms of F1 score obtained across 5 independent runs is reported.

From the Table, we see that the best performance are obtained by Xception and MobileNet, the two popular deep learning architectures using separable 2D convolutions. Such a result, encouraged us to adopt in our experiments an architecture with residual connections and SepConv2D layers, similar to Xception and MobileNet.

II. COMPARISON OF DIFFERENT TECHNIQUES TO HANDLE CLASS IMBALANCE

Tab. V reports the classification performance and training times when using different techniques to handle class imbalance. Despite being more computationally intensive, the oversampling technique yields the best classification performance and, thus, is the one adopted in the experimental evaluation.

Balancing method	Time/epoch	TN	FN	FP	TP	F1 score
Class-weighting	18s	369.4 \pm 2.1	12.3 \pm 1.6	2.3 \pm 1.3	52.1 \pm 1.9	0.88 \pm 0.01
Oversampling	31s	367.8 \pm 2.7	3.2 \pm 2.7	4.6 \pm 1.7	59.4 \pm 1.7	0.94 \pm 0.01
Rejection sampling	22s	367.3 \pm 2.4	6.1 \pm 1.4	4.2 \pm 1.1	58.0 \pm 2.1	0.92 \pm 0.02

TABLE V: Classification performance obtained with different methods to handle class imbalance.

III. OPTIMAL HYPERPARAMETERS FOR MODELS TRAINED ON LOWER RESOLUTION IMAGES

We optimized the hyperparameters for the models trained on lower resolution by following the exact same procedure that we used for the model operating on the higher resolution images. The optimal hyperparameters for the different models are reported in the Tab. VI. We note that the optimal hyperparameters are the same for different image resolutions, except for: *Num. residual blocks*, *Dropout rate*, and *Learning rate*.

Hyperparam.	500m	1km	2km
Activation	ReLU	ReLU	ReLU
Conv2D filters	8	8	8
Kernel size	3	3	3
SepConv2D filters	8	8	8
Num. residual blocks	7	5	4
Global pooling	avg	avg	avg
Num. dense layers	1	1	1
Units in the dense layer	8	8	8
Dropout rate	0.5	0.4	0.6
Use batch normalization	True	True	True
Learning rate	1e-3	1e-2	1e-3

TABLE VI: Optimal hyperparameters for the models trained on different image resolutions.

IV. COMPARISON BETWEEN CO- AND CROSS-POLARISED CHANNELS

Here, results using only one polarisation channel are presented for the 500 m resolution images. Specifically, the F1-score is presented in table VII (averaged over 5 independent runs), obtained when using the co- or cross-polarised channels separately. For the co-polarised case, both HH and VV are used jointly, while for the cross-polarised case, the VH and HV channels are used jointly. The results show that the cross-polarised channels yields better performance compared to the co-polarised channels. A possible reason for this could be that the cross-polarised channels better captures high wind speed features. At

¹²<https://keras.io/api/applications/>

high wind speeds, the co-polarised backscatter saturates faster as a function of wind speed [57]. Another factor that could play a role is that the two co-polarisations (HH and VV) typically behaves differently as a function of incidence angle and target properties, making the dataset somewhat heterogeneous. The cross-polarisation (VH and HV) dataset is on the other hand homogeneous, since, theoretically these channels are identical for reciprocal targets, like the ocean.

The best performance is however still obtained using both polarisation channels, as shown in table II in the main manuscript.

	Co-pol ($x_{ }$)	Cross-pol (x_{\times})
F1-score	0.886 ± 0.016	0.916 ± 0.013

TABLE VII: Results obtained using only co- or cross-polarised channels separately.

V. TRAINING TIME FOR DIFFERENT IMAGE RESOLUTIONS

Tab. VIII reports the training times of the proposed deep learning architecture when images of different resolutions are used in training. The training times are measured on an Nvidia RTX 3090. Clearly, lower resolutions result in a much faster training. However, even when using 500m resolution, the neural network can be trained reasonably fast.

About the differences in time for the inference phase, they are negligible when using different image resolutions (a fraction of a second in each case). Considering the whole process from satellite acquisition, data downlink/download, SAR focusing, pre-processing (in particular geocoding) etc, the inference time of the neural network model is by all means negligible in an operational setting.

Pixel size	500m	1km	2km
Time/epoch	31s	9s	5s

TABLE VIII: Training times for different image resolutions.

VI. INTERPRETABILITY FOR A MODEL TRAINED ON LOW-RESOLUTION IMAGES

An interesting question when comparing model performance between input image resolutions (500, 1000 and 2000 metres), is if the interpretability metrics (IG and Grad-CAM) can indicate why the performance is worse for the lower resolutions. By comparing the results of the high-res model trained on the 500m resolution images to the results of the low-res model trained on the 2000m resolution images, we find that the low-res model miss-classifies 11 FNs and 7 FPs. Among these, there are 9 FNs and 5 FPs that the high-res model classifies correctly. Therefore, we compare the low-res FNs to corresponding high-res TPs, and low-res FPs to high-res TNs.

Since we cannot compute IG on samples classified as negatives (as explained in section IV-C), we only consider the Grad-CAM. It should however be noted that the Grad-CAM heat map is a result of gradients at the last layer in the model (see section III-E2 for details). Since each layer in the model contains pooling, the heat map will be of lower resolution than the input image itself. The Grad-CAM heat map thus provide little or no information about fine detailed differences. Yet, it is expected that differences between the high- and low-res results are primarily fine details (which disappear when the resolution is lowered). Despite this limitation in the analysis, differences in the interpretability results with regard to Grad-CAM are presented below.

A. Low-res FNs versus high-res TPs

The FNs of the low-res model (that the high-res model classifies correctly), could give insights into what key features of the input images are lacking at the lower resolution in order to correctly classify an image with a mesocyclone. The low-res FNs and the corresponding high-res TPs are shown in figure 12 (3 of the total of 9 cases are shown). It is clear that the low-res model does not attribute the same importance to the cyclonic eye or wind front features as the high-res model. This could indicate that at the cyclonic eye or at the wind front, high-res features are of particular importance. If these are lacking, the model focuses elsewhere in the image. In the shown examples, the attribution of the low-res model appears considerably more scattered, which could indicate that there are not sufficiently strong features to attract the model attention.

B. Low-res FPs versus high-res TNs

The low-res FPs and the corresponding high-res TNs are shown in figure 13 (3 of the total of 5 cases are shown). It is clear that the low-res model now puts the main attribution to the image centres, while the attribution of the high-res model is more scattered. In the top sample, part of the image is covered by sea ice, which is rich of fine details. At the lower resolution, these features could potentially be sufficiently blurred for the low-res model to be confused, e.g., with a wind front. In the second sample, a slight wind front seems to be picked up, but it is unclear what the low-res model actually is focusing on. The situation is similar in the third example, where no clear feature is shown at the centre.

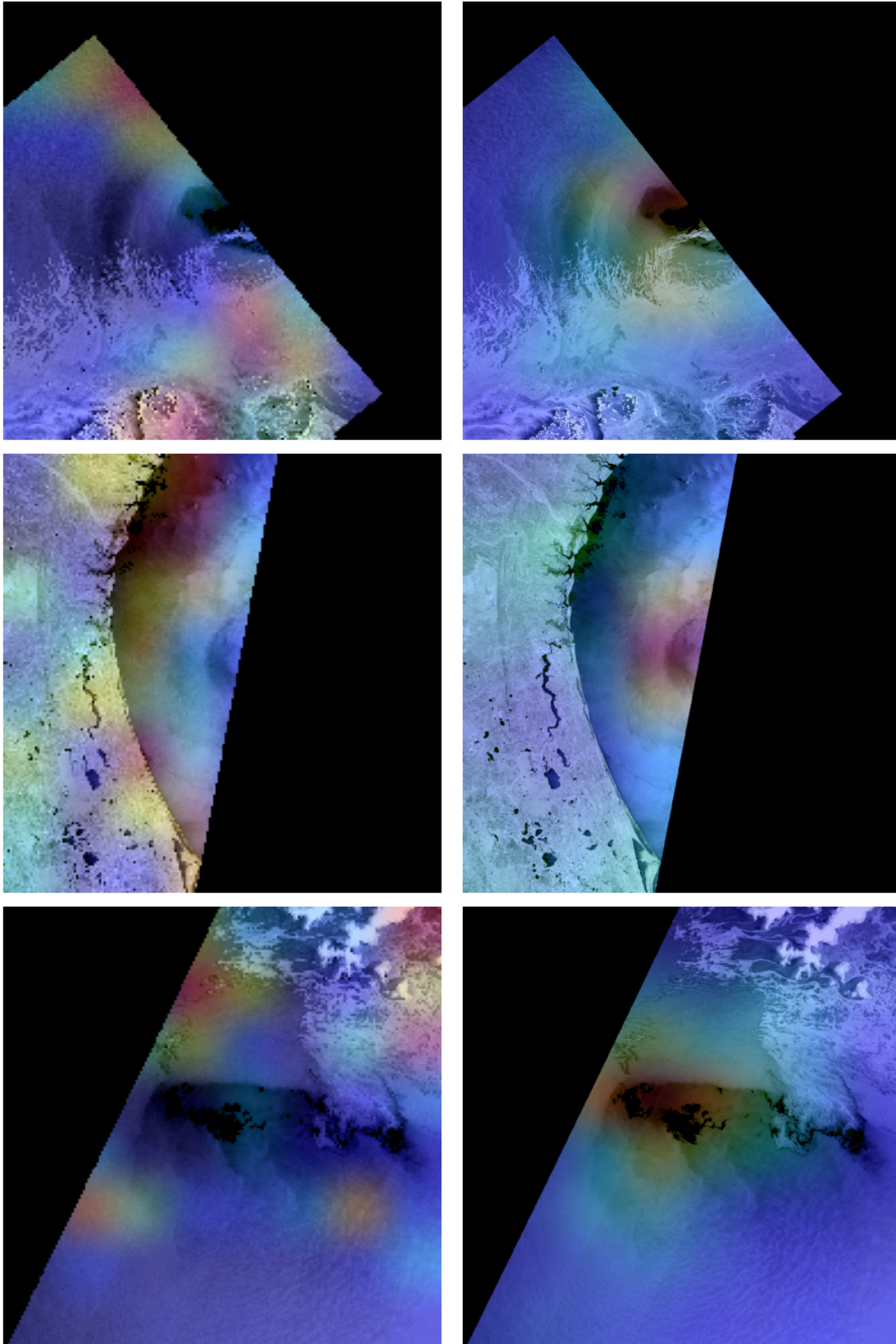


Fig. 12: Grad-CAM heat maps for the FNs of the low-res model to the left, and the corresponding TP of the high-res model to the right.

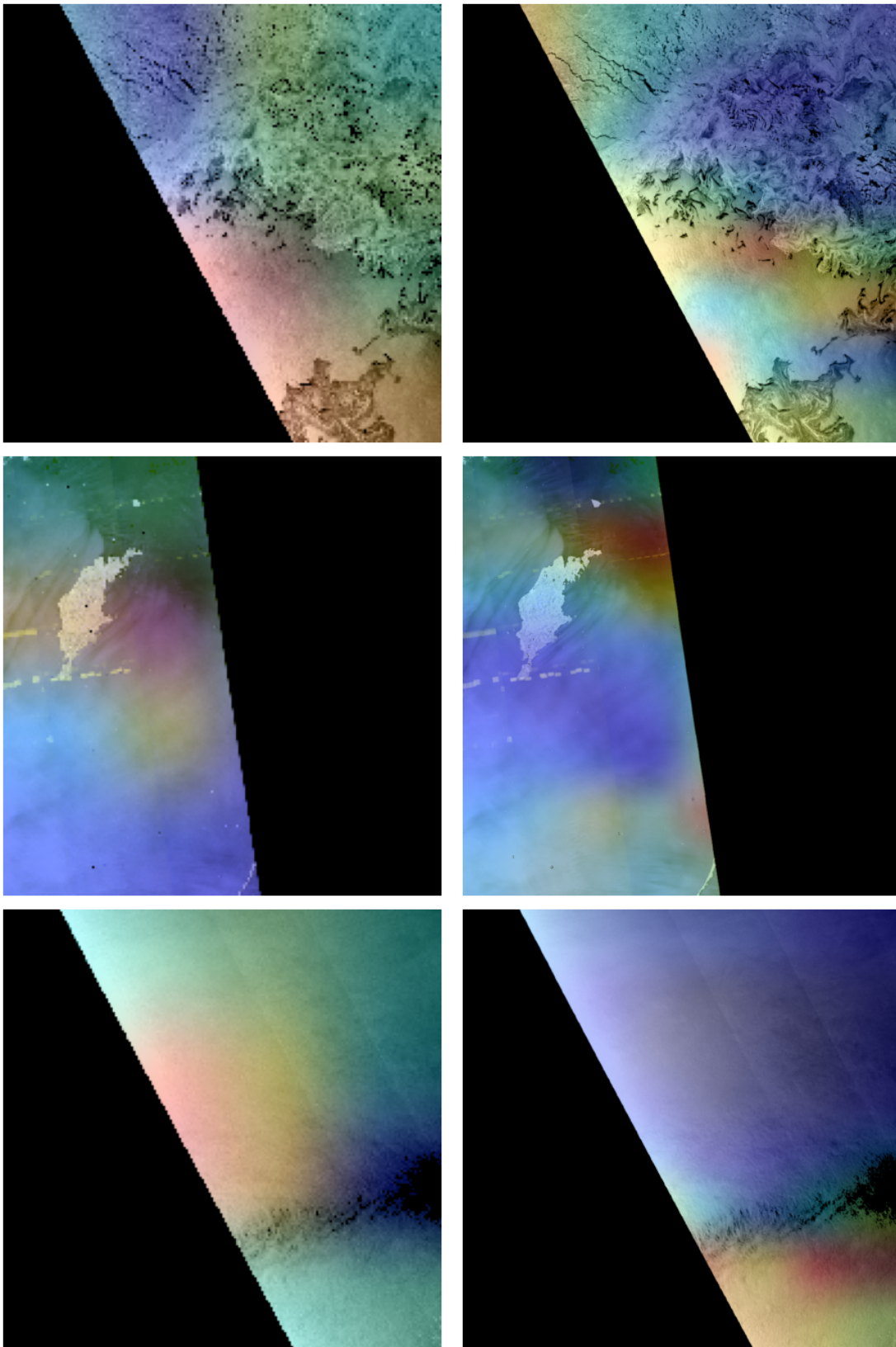


Fig. 13: Grad-CAM heat maps for the FPs of the low-res model to the left, and the corresponding TN of the high-res model to the right.