

1 Marine cloud base height retrieval from MODIS cloud properties using 2 machine learning

3

4 Julien LENHARDT ¹, Johannes QUAAS ^{1,2}, Dino SEJDINOVIC ³

5

6 ¹Leipzig Institute for Meteorology, Leipzig University, Leipzig, Germany

7 ²ScaDS.AI - Center for Scalable Data Analytics and Artificial Intelligence, Leipzig University, Humboldtstraße 25, 04105

8 Leipzig, Germany

9 ³School of Computer and Mathematical Sciences & Australian Institute for Machine Learning, University of Adelaide, Adelaide,

10 Australia

11 *Correspondence to:* Julien LENHARDT (julien.lenhardt@uni-leipzig.de)

12 Abstract

13

14 Clouds are a crucial regulator in the Earth's energy budget through their radiative properties, both at the top-of-the-atmosphere
15 and at the surface, hence determining key factors like their vertical extent is of essential interest. While the cloud top height is
16 commonly retrieved by satellites, the cloud base height is difficult to estimate from satellite remote sensing data. Here we present
17 a novel method called ORABase (Ordinal Regression Autoencoding of cloud Base) leveraging spatially resolved cloud
18 properties from the MODIS instrument to retrieve the cloud base height over marine areas. A machine learning model is built
19 with two components to facilitate the cloud base height retrieval: the first component is an autoencoder designed to learn a
20 representation of the data cubes of cloud properties and reduce their dimensionality. The second component is developed for
21 predicting the cloud base using ground-based ceilometer observations from the lower dimensional encodings generated by the
22 aforementioned autoencoder. The method is then evaluated based on a collection of co-located surface ceilometer observations
23 and retrievals from the CALIOP satellite lidar. The statistical model performs well on both datasets, exhibiting accurate
24 predictions in particular for lower cloud bases and a narrow distribution of the absolute error, namely 379 m and 328 m for the
25 mean absolute error and the standard deviation of the absolute error respectively for cloud bases in the test set. Furthermore,
26 cloud base height predictions are generated for an entire year over ocean, and global mean aggregates are also presented,
27 providing insights about global cloud base height distribution and offering a valuable dataset for extensive studies requiring
28 global cloud base height retrievals. The global cloud base height dataset and the presented models constituting ORABase are
29 available from Zenodo (Lenhardt et al., 2024).

30 1 Introduction

31

32 Clouds play a key role in the Earth's energy budget through their interactions with incoming shortwave and outgoing longwave
33 radiation fluxes. It is thus critical to adequately quantify cloud radiative properties and their changes under global climate
34 change. However, cloud radiative properties remain a large uncertainty in estimating anthropogenic climate change and possible
35 impacts in the future (Boucher et al., 2013; Forster et al. 2021). Radiative properties of clouds are related to numerous quantities
36 that can be used to characterise them. For instance, the cloud base height (CBH) is a crucial radiative property through its impact
37 on the surface longwave radiation. Furthermore, the cloud geometrical thickness (CGT), defined as the difference between the
38 cloud top height (CTH) and the CBH, links to the adiabatic cloud water content allowing the quantification of the cloud's
39 subadiabaticity. Additionally, deriving the CBH is of practical use for pilots, providing crucial information during flights.

40 However, while the CTH can be rather easily obtained through passive satellite observations, the CBH retrieval remains
41 problematic due to the fact that it is only indirectly accessible to satellites, and due to retrieval errors related to satellite remote
42 sensing such as instrument shortcomings or noisy measurements. Since the difference between the CTH and the CBH quantifies
43 the vertical extent of a cloud, one way to retrieve the CBH from passive satellites is by making heavy assumptions on the vertical
44 distribution of the cloud water path inside the cloud profile. It is thus a challenging retrieval with passive satellites data that
45 provide information about the cloud top (e.g. cloud top temperature (CTT), pressure (CTP) or height (CTH)) or about the entire
46 column (e.g. cloud optical thickness (COT)) assuming the cloud's adiabaticity. For example, Noh et al. (2017) rely on a
47 semiempirical approach to link the CGT to the CTH and the cloud water path (CWP, includes both ice and liquid water paths). In
48 a different approach, Böhm et al. (2019) retrieve the CBH from triangulation of a multi-angle spectroradiometer. However, in
49 this case, assumptions were required on the distribution of convective clouds. On the other hand, active satellite remote sensing
50 retrieves information with vertical resolution which greatly helps resolving the clouds vertical distribution. However, active
51 satellite measurements can display attenuated signals close to the surface (Tanelli et al., 2008; Marchand et al., 2008) particularly
52 in the presence of thick clouds or precipitation, rendering the retrieval of the CBH difficult even for radar and lidar. Among
53 others, Mülmenstädt et al. (2018) and Lu et al. (2021) present methods focusing on low clouds which use the CBH from active
54 satellite retrievals of neighbouring thin clouds as representative of the surrounding cloud field. Active remote sensing
55 additionally suffers from the sparse sampling that is confined to a narrow swath below the satellite. Finally, Goren et al. (2018)
56 combine information from both passive and active satellite remote sensing and rely upon an adiabatic cloud model to derive the
57 CBH. The retrieval of the CBH using satellite remote sensing data relies on a number of simplifying assumptions and is,
58 consequently, prone to errors. Subsequently, uncertainties in the estimation of the CBH propagate into uncertainties in the overall
59 cloud radiative effect (CRE) (Kato et al., 2011; Trenberth et al., 2009).

60 The method presented here called ORABase (Ordinal Regression Autoencoding of cloud Base) leverages passive satellite
61 retrievals of cloud properties in combination with marine surface observations to derive the CBH of a cloud scene using a
62 machine learning (ML) model. The CBH retrieval method relies on level 2 satellite data, namely three different cloud properties
63 which are CTH, COT and CWP. A convolutional neural network (CNN, LeCun et al., 1989; LeCun et al., 1995) model following
64 the autoencoder (AE; Kramer, 1991; Hinton et al., 2006) framework is trained in a self supervised way to reconstruct the
65 previously mentioned cloud properties. This type of artificial neural network has been widely used in computer vision
66 (Krizhevsky et al., 2012; LeCun et al., 2010) but also more recently in various applications in climate science (Reichstein et al.,
67 2019; Watson-Parris et al., 2022). Thereafter, an ordinal regression (OR; Winship et al., 1984) model is fitted to predict the CBH
68 corresponding to the cloud properties, learning from ground-based marine CBH retrievals. These different steps constituting the
69 method are summarised in [Figure 1](#) and detailed in section 2. The objective of the developed method is primarily to produce
70 CBH retrievals with reduced uncertainty, and additionally to extrapolate CBH retrievals from local surface observations to a
71 wider spatial and temporal coverage. Indeed, we hypothesise that the spatial pattern of the cloud field carries information about
72 the CBH and that the CNN can exploit the potential non-linear relationship between the CBH and the satellite observations.
73 Furthermore, as more accurate CBH retrievals are obtained from ground-based remote sensing observations which are only
74 available at isolated locations, we capitalise on these retrievals to develop a satellite-based retrieval algorithm capable of
75 generalising to global distributions. We sensibly reduce the scope of the study by focusing on lower clouds, in particular as
76 ground-based CBH observations display higher accuracy compared to satellite-based retrievals in those cases, and as it is the
77 lowest cloud which often matters most for e.g. the surface radiation budget. We also restrict the retrievals to marine regions to
78 remove the impact of orography on surface observations especially for these same low level clouds.

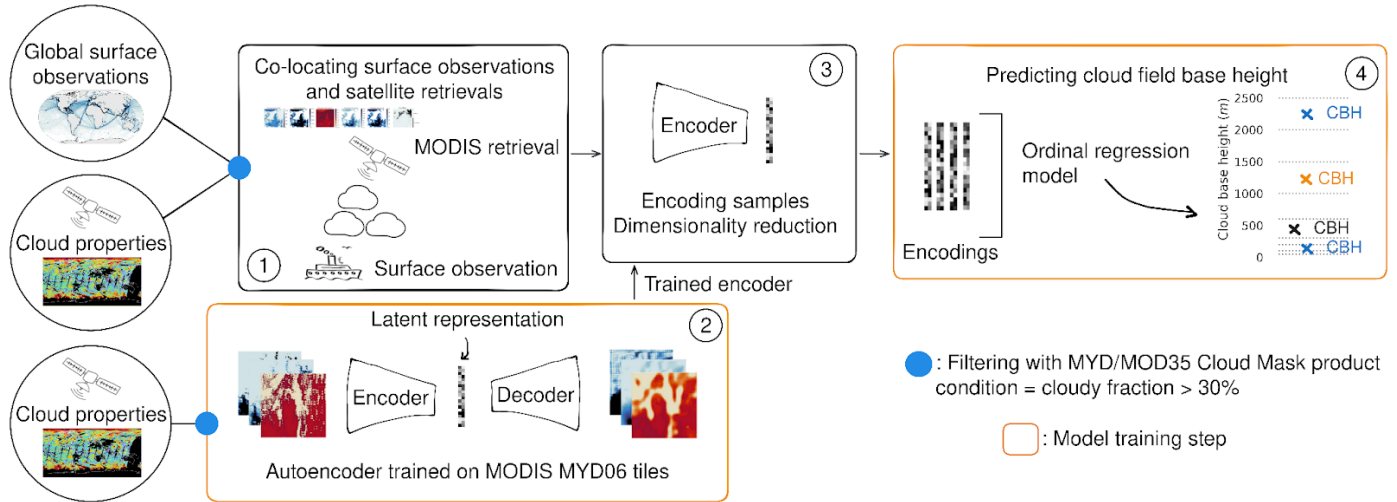
79 Section 2 firstly introduces the datasets and the co-location between ground-based observations and satellite retrievals. Secondly,
80 the ML method constituting ORABase is described. In section 3 we evaluate our predictions against other methods including
81 Noh et al. (2017) and other products from active satellite measurements like the 2B-CLDCLASS-LIDAR product (Sassen et al.,

82 2008). Section 4 presents the global dataset of the CBH which is derived from the ML approach. We discuss the benefits and
 83 remaining challenges of our method in section 5. Further details about the spatial distribution of the observations and the ML
 84 method are included in the appendices A-E. Additional links to available data outputs and codes are listed in the corresponding
 85 sections.

86

87 2 Data and methods

88



89

90 **Figure 1: Schematic of the cloud base height retrieval method. 1) Co-location of surface-based cloud base height**
 91 **observations and satellite retrievals. 2) Autoencoder training on satellite cloud properties. 3) Encoding of co-located**
 92 **samples using the trained encoder. 4) Prediction of the cloud field base height.**

93

94 2.1 Surface observations

95

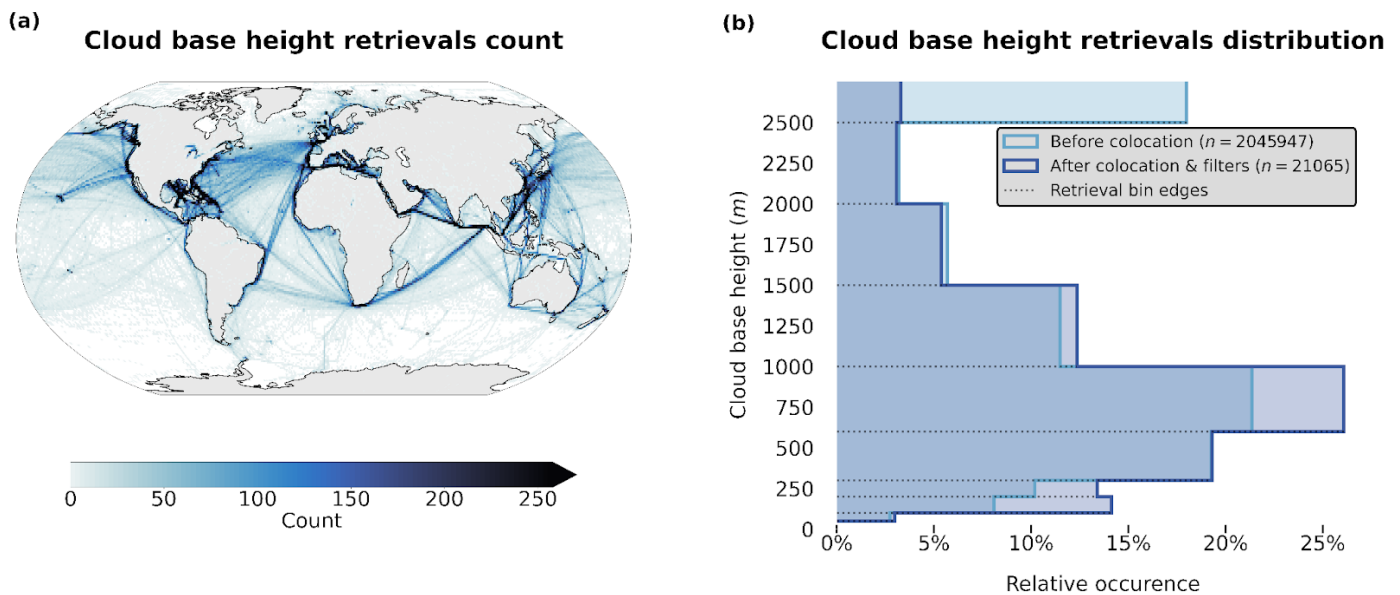
96 The CBH labels used in this study are part of a global marine meteorological observation dataset maintained by the UK Met
 97 Office (Met Office, 2006; [Table 1](#)), which provides observational data ongoing from 1854. The observations are conducted from
 98 measuring stations that were located on ships, buoys or platforms. As a consequence, this study largely relies on observational
 99 data representing the areas along the corresponding ship routes ([Fig. 2a](#)). Despite their coarse resolution, the reported cloud base
 100 observations provide valuable information about clouds in remote marine areas. The distribution of CBH observations and
 101 corresponding bins are shown in [Figure 2](#).

102 At the beginning of meteorological and weather reports, surface-based cloud observations were retrieved manually or visually by
 103 human observers, but they have been gradually replaced by automated systems. The CBH is derived using a ceilometer, an
 104 instrument based on a laser pointing upright and measuring the backscatter from the cloud base, and is then reported following
 105 the current standards from the World Meteorological Organisation (WMO; WMO, 2019). The CBH observations are sorted into
 106 bins of increasing width (from 50 m to 500 m bin width) corresponding to the altitude ([Fig. 2b](#)) as the data transfer through radio
 107 limits the amount of transferable information and precision close to the surface is of importance notably for aircrafts. Since the
 108 actual measured CBH values are not available in the dataset, it is impossible to directly quantify a possible bias stemming from
 109 this binning process. In general here, we can suspect that the available CBH retrievals represent an accurate or underestimated
 110 assessment of the effective CBH, as for example a ceilometer measuring a CBH of 2490 m will be reported in the 2000 m bin in
 111 the available dataset. Using for example the central value of each bin could be another way to compute averages to potentially
 112 alleviate this unknown bias but it is not presented here. However, the method presented in the following sections predicts the
 113 CBH in corresponding bins, so it is left to the user to use these as they see fit for further analysis.

114

Data product	Description	Variables	Resolution	Usage
Global marine meteorological observations (Met Office, 2006)	Surface observations	Cloud base height (m)	Latitude/longitude coordinates 0.1° Hourly/daily observations	Labels
MODIS Atmosphere L2 Cloud Product (MYD06) (Platnick et al., 2017)	Cloud-top properties, cloud optical and microphysical properties	Cloud top height, CTH (m) Cloud optical thickness, COT (a.u.) Cloud water path, CWP (g.m ⁻²)	1 km pixel resolution Daily overpass	Input features
MODIS Atmosphere L2 Cloud Mask Product (MYD35) (Ackerman et al., 2017)	Cloud pixel flag	Cloud mask	1 km pixel resolution Daily overpass	Used for cloud scene filtering

116 **Table 1 : Dataset description. The MODIS data are derived from the collection 6.1 of the datasets (Platnick et al., 2017;**
 117 **Ackerman et al., 2017; cf. section 2.1). The surface observations are provided by a worldwide station network available**
 118 **from the UK MetOffice (Met Office, 2006; cf. section 2.2).**
 119



120

121

122 **Figure 2: (a) Spatial distribution of cloud base retrievals count (1° grid) and (b) distribution of the retrieved cloud base**
 123 **height before and after the co-location and filtering process, for observations from the years 2008 and 2016.**

124

125 2.2 Satellite data

126

127 In this study we use products from the MODerate Resolution Imaging Spectroradiometer (MODIS, Platnick et al., 2017) from
 128 the AQUA satellite as input data that is later combined with the CBH labels derived from the surface-based observations to train
 129 the prediction model. We choose MODIS satellite retrievals as they provide a large amount of data with kilometre-scale
 130 resolution and daily overpasses, the spatial coverage of one granule representing an area of 2330 km x 2000 km. We make use of
 131 the CUMULO dataset (Zantedeschi et al., 2019) since it provides already preprocessed satellite data from the A-train with daily
 132 full coverage of the Earth for the years 2008 and 2016. In particular out of the available variables we use two aligned products
 133 (cf. [Table 1](#)), namely the MODIS06 level 2 cloud product (hereafter MYD06; Platnick et al., 2017) which provides relevant

134 cloud properties and the MODIS35 level 2 cloud flag mask (hereafter MYD35; Ackerman et al., 2017) which allows us to filter
135 scenes and screen for clouds.

136 The MYD06 product contains various cloud top properties (temperature, pressure, height) and cloud optical and microphysical
137 properties (optical thickness, effective radius, water path). Level 2 data are derived from calibrated radiances through various
138 algorithms and physical relations detailed in Platnick et al. (2017). The cloud top quantities are derived from radiance data of
139 several channels. Wavelengths in the CO₂ absorption range are particularly used to identify the cloud top pressure (CTP) and thus
140 the CTH of high clouds because of the opacity of CO₂. For thicker or low boundary layer clouds, since the CO₂ slicing technique
141 fails, the CTH is retrieved using the 11 μm brightness temperature band and combined with simulated brightness temperatures
142 based on vertical profiles from GDAS using surface temperature together with monthly averaged lapse rate data (Baum et al.,
143 2012). The use of monthly averaged lapse rate data separately for different regions greatly helped reduce the bias in retrieved
144 CTHs for low clouds in the Collection 6 of MYD06 from Collection 5, but some spatial and regional biases remain. These biases
145 directly impact the spatial and temporal distribution of CTH in the data and thus what the model could learn from. The cloud
146 optical thickness (COT) and cloud effective radius (CER) are simultaneously derived from multispectral reflectances, cloud
147 masks, CTP data and surface type characteristics. The cloud water path (CWP) is additionally retrieved as part of the cloud
148 optical properties algorithm described in Platnick et al. (2017). The retrieval of these cloud properties additionally requires inputs
149 such as temperature, water vapour and ozone profiles from NCEP GDAS (Platnick et al., 2003; Baum et al., 2012) which can
150 lead to potential uncertainties in particular in remote marine regions where only sparse observations are available for
151 assimilation.

152 In general, the MYD06 level 2 product offers the advantage that the statistical model can be built relying on cloud properties and
153 it can thus allow the study of relationships between the CBH and other cloud properties. Calibrated radiances, one step ahead in
154 the data processing pipeline, would also provide insightful information but would require inputs of larger dimensionality since
155 key information about clouds would be scarcer. Furthermore, using MYD06 level 2 data allows us to compare our method to
156 others which in most cases use cloud properties to retrieve the CBH. The level 2 product provides pre-processed data on top of
157 the calibrated radiances and reflectances of level 1 data, which might introduce biases in the statistical model as previously
158 mentioned regarding the CTH for example. From the entirety of available MYD06 retrievals, we select three cloud properties in
159 particular, namely the CTH, COT, and CWP. The CTH is used as it provides key information about the CBH in the cloud field, as
160 seen in Böhm et al. (2019). Vertically integrated cloud quantities like the COT and CWP further help the statistical model by
161 providing key information about the cloud's vertical extent, lacking in cloud top only properties, making them commonly used
162 for retrieving the CBH (e.g. Noh et al., 2017). The CWP as computed from COT and CER, and, in consequence, also the CBH
163 are built on adiabatic assumptions (Grosvenor et al., 2018) and therefore cannot be used to constrain subadiabaticity as also
164 highlighted in Mülmenstädt et al. (2018).

165

166 2.3 Datasets co-location

167

168 We proceed to collocate our two data sources over the two years of MODIS MYD06 data available. To obtain the cloud properties
169 of the cloud scene corresponding to the surface retrieval of CBH, we select a square tile of 128 km x 128 km from the *closest*
170 MODIS granule available centred around the observation location. Here *closest* means that the MODIS granule contains the
171 (latitude, longitude) coordinate of the CBH observation and the full extent of the tile centred around, and that the satellite
172 retrieval was made during a one hour time-window before/after the CBH observation time. The spatial and temporal thresholds
173 used to collocate the surface observations and the satellite retrievals are chosen for several reasons. Mainly, we want the satellite
174 cloud properties to be representative of the cloud scene for which the CBH observation was made. Additionally, we want to
175 recover a satisfying number of samples during the collocation process. Further arguments regarding the sensitivity of the retrieval
176 method to the tile size are described in the following method section 2.5.

177 The extracted tile corresponding to the surface observation is then filtered. A first filter is applied to missing values in the
178 different cloud properties fields to primarily avoid retrievals of poor quality. This is predominantly the case for the COT and
179 CWP fields for which the retrieval fails more frequently, sometimes entirely. Another filtering is concordantly done using the
180 MYD35 product for cloud cover (minimum of 30% of cloudy pixels) to ensure the cloud field was substantial enough for the
181 collocated surface observation to be representative. Additional comments on the sensitivity of the CBH retrieval to this threshold
182 are presented in the following section on the downstream task of CBH prediction. Throughout the quality filtering process, the
183 missing data is one of the major factors impacting the amount of retained samples. On [Figure 2](#), we can see that it seems to
184 impact the clouds with higher CBHs.

185 The overall filtering and co-location process yields around 21 000 samples. This only represents around 1% of the initial CBH
 186 observations mainly due to the co-location process both in time and space with the MODIS overpasses. Missing values and cloud
 187 cover filters are an additional factor in the reduced number of co-located samples. The presented co-located dataset is the basis to
 188 build our cloud scene CBH retrieval.

189

190 2.4 Autoencoder

191

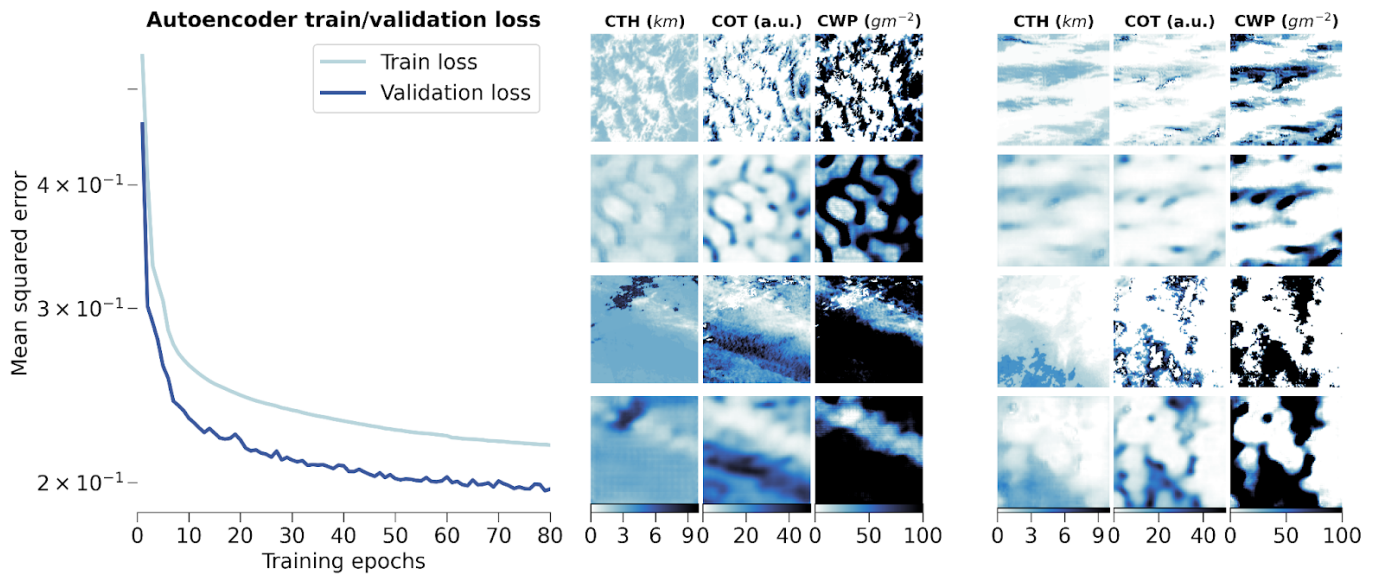
192 To circumvent the lack of labelled samples from which the relevant features are extracted, and to learn useful lower-dimensional
 193 representations of the data, we add a dimensionality reduction step to our method through an unsupervised learning model. AEs
 194 offer a wide application spectrum, ranging from preprocessing to the generation of new outputs. AEs are commonly used in
 195 unsupervised learning settings for reducing the dimension of the input data to leverage the latent representations learned by the
 196 model to perform clustering, classification or regression in a lower dimensional space (Baldi et al., 2012). We use classical AEs
 197 for their simplicity and versatility, but other approaches to unsupervised latent representation learning, such as variational AEs
 198 and its many variants, can be used in a similar fashion. In general, AEs learn to encode the given input data to produce a latent
 199 representation of lower dimension. From the latent representation, the input data is then reconstructed. The learning process is
 200 driven by what is called the reconstruction loss that minimises the difference between the input and the reconstructed output.

201 Here we use a convolutional AE architecture which is based on a CNN backbone in order to leverage the spatial structure of our
 202 input data (Pu et al., 2016). In particular, we rely on the widely employed CNN architectures U-Net (Ronneberger et al., 2015)
 203 and VGG (Simonyan and Zisserman, 2015), where the convolution layers are based on 3x3 filters, stacked in blocks followed by
 204 maximum pooling layers, and mirrored for the decoder part of the model using transposed convolution layers (Zeiler et al.,
 205 2010). We adapt the size of the input to fit our chosen tile size (128), the latent space size to 256, and use the improved Leaky
 206 Rectified Linear Units (LeakyReLU; Maas et al., 2013) over the original ReLU (Nair and Hinton, 2010) as activation functions.
 207 The detailed parameterization of the model is described in Appendix B. The model code was developed following
 208 implementations from the packages *PyTorch* (Paszke et al., 2019) and *TorchVision* (TorchVision, 2016) and is included in the
 209 related Zenodo archive (Lenhardt et al., 2024). The main goal of the AE training is then to minimise the loss function during the
 210 optimization or learning process, and to reproduce the input data with the highest fidelity. For the loss function which in this case
 211 is only the reconstruction error, we use the common mean-squared error (MSE), which can be written for a batch of samples as :

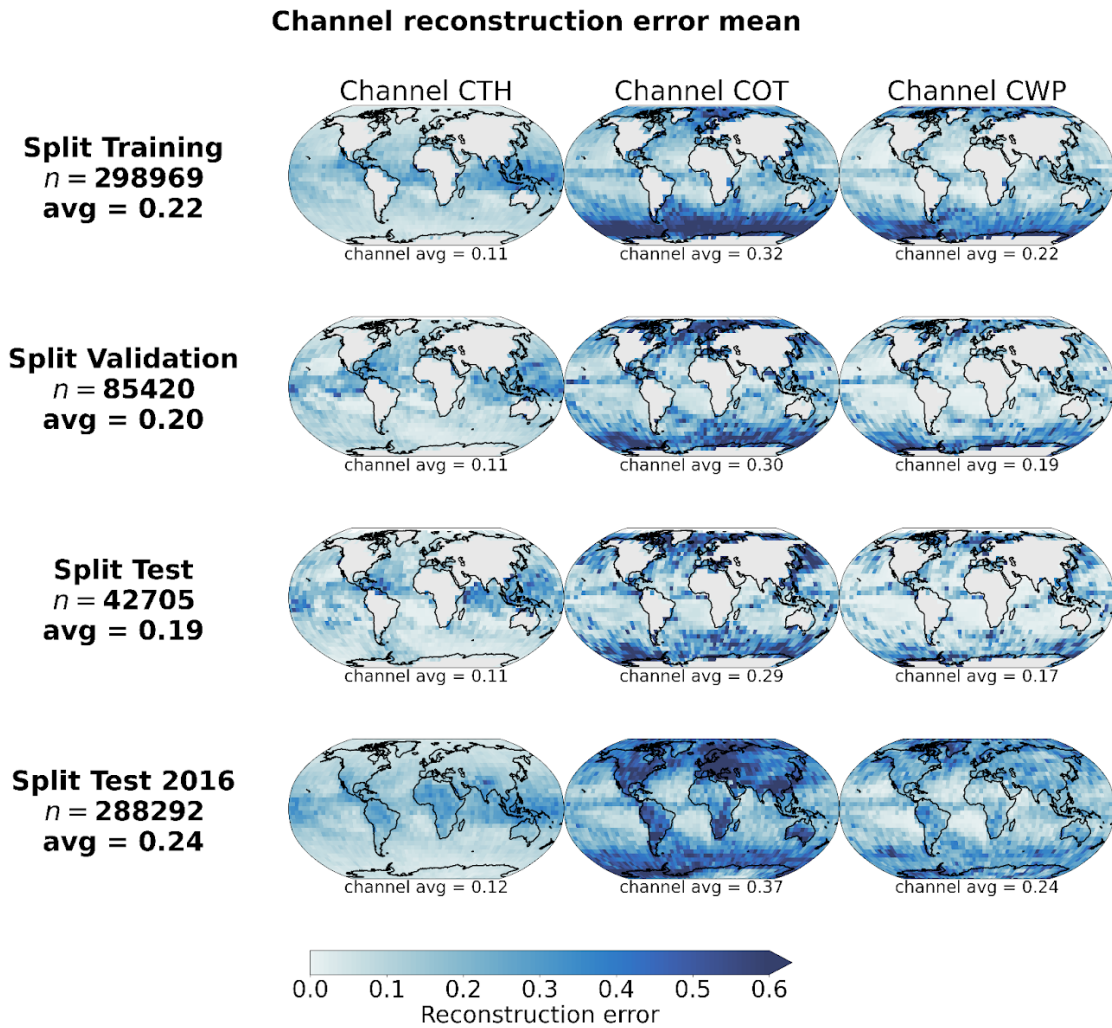
$$212 \quad \mathcal{L}_{reconstruction} = \frac{1}{N_i} \sum_{b \in B_i} \left\| b - D_{\theta}(E_{\theta}(b)) \right\|_2^2 \quad (1)$$

213 where, with the tiles used for training the AE noted as $B = \{b_n \in \mathbb{R}^{3 \times 128 \times 128}\}_{n \in [1, N]}$, B_i represents a batch of samples of size
 214 N_i and θ the combined parameters of the encoder E and decoder D models.

215 However, this self supervised step requires a large amount of data that the AE can learn from. Therefore, we select one full year
 216 of data of MODIS granules from the CUMULO dataset (from the year 2008, cf. section 2.2) and randomly sample tiles following
 217 the same criteria as during the co-location process (cf. section 2.3). We sample a maximum of 20 tiles from a single granule and
 218 this for only a single year of data in order to avoid possible spatial and temporal auto-correlation in the data used for training and
 219 testing leading to a non-representative performance of the mode (Kattenborn et al., 2022). Further details on the study of the
 220 generalisation performance of the model for new observations in space and time are given in appendix B. The overall built
 221 dataset consists of around 500 000 samples which are then splitted for training, validation and testing based on their retrieval
 222 date. We additionally create a dataset based solely on data from the year 2016 for further testing which includes tiles not only
 223 over ocean but also over land, indicating potential generalisation skill for unseen data including orography influence. The
 224 reconstruction error during training and validation is shown in [Figure 3](#) along with examples of reconstructed samples. The
 225 spatially averaged reconstruction errors per cloud property channel are displayed in [Figure 4](#) for each of the training, validation
 226 and testing datasets previously mentioned. The trained model reaches an MSE of 0.19 on the test set and of 0.24 on the global
 227 test set of 2016. The presented model is trained on tiles of size 128x128, but some arguments regarding the choice of the tile size
 228 are made in the following section in the context of the downstream task of CBH prediction.



229
 230 **Figure 3: (left) Training and validation losses during model optimization. (right) Examples of tiles (first and third rows)**
 231 **with the corresponding reconstructions (second and fourth rows) for the different cloud property channels.**
 232



233
 234 **Figure 4: Spatial distribution of channel reconstruction errors aggregated on a 5° grid for the 2008 training, validation,**
 235 **test and the 2016 test datasets.**
 236

237 2.5 Cloud base height ordinal regression

238

239 Once the AE's optimization process is completed, the next step is to predict the corresponding CBH for the observed scene. As
240 seen in [Figure 2](#), the retrieved CBH observations are binned into different categories following WMO standards (WMO, 2019).
241 This leads to a prediction problem at the intersection of regression (i.e. predicting numerical values) and classification (i.e.
242 predicting the object class) called ordinal regression (OR). The labels from the target variable are defined by classes following a
243 certain order, in this case the increasing CBH. A wide array of methods stems from this field with diverse applications for
244 example in computer vision using neural networks (e.g. Niu et al., 2016; Shi et al., 2023; Lazaro and Figueiras-Vidal, 2023).
245 Different methods exist to tackle such problem setups either via modification of the target variable, ordinal binary decomposition
246 or threshold modelisation (Gutiérrez et al., 2016; Pedregosa et al., 2017). Threshold models were shown to be able to perform
247 better than the ones designed for regression or multi-class classification on OR tasks (Rennie et al., 2005). We consider here two
248 alternative frameworks in the case of threshold models which differ in how they penalise threshold violations:
249 immediate-threshold (IT; [Eq D.1](#)) and all-threshold (AT; [Eq D.2](#)). The overall training process of the model aims at optimising a
250 set of weights to project the input data to a one dimensional plane, subsequently dividing the constructed representation using
251 learnable thresholds. These two implementations of threshold models are available from the *mord* Python package (based on
252 Pedregosa, 2015) and further details on threshold OR models are added in appendix D.

253 To help evaluate the prediction model, we rely on a set of different metrics pertaining either to the regression aspect of the
254 problem or to its classification/ordinal nature. First, the macro-averaged mean absolute error (MA-MAE) is used as it weights
255 each class separately before averaging the subset MAEs, making it useful in the case of OR problems with imbalanced datasets
256 (Bacianella et al., 2009). Using a macro-averaged metric prevents us from choosing a trivial model which might always predict
257 the dominating class. Additionally, the macro-averaged root mean square error (MA-RMSE) is also used to investigate the skill
258 of the prediction models. To assess the ordering of the predicted retrievals with respect to the labels, the ordinal classification
259 index (OC; Cardoso and Sousa, 2011) and its updated version the uniform ordinal classification index (UOC; Silva et al., 2018)
260 are computed. A version of the latter not requiring an extra hyperparameter, the area under the UOC (AUOC; Silva et al., 2018),
261 is also reported. These different metrics are able to capture the proper ranking order of the predictions compared to the labels
262 using the confusion matrix and also the overall accuracy of the prediction model. Nevertheless, one caveat is that these indexes
263 developed for ordinal classification assume each class to be equally distant from another which is not the case here since the
264 CBH retrievals are reported in bins of variable width. However, a purely ordinal classification index will drop all information on
265 the scale of the response (1500 m misclassified as 600 m treated the same as 200 m misclassified as 50 m, since only the order
266 matters) which might be not entirely appropriate for this problem. In an effort to address this limitation, the indexes are adapted
267 to mimic the spacing between the different CBH bin classes by incorporating classes that are all spaced by 50 m, ranging from 50
268 m up to 2500 m. In this manner, the CBH class difference is more suited to the actual nature of the retrieval.

269 However, several aspects of the ordinal regression model need to be investigated first. To this extent, we first divide our global
270 colocated dataset (section 2.3) in training, validation and testing datasets but while ensuring each class is relatively equally
271 represented in each split. The following aspects and sensitivities of the model to the input data parameters are assessed using the
272 training and validation datasets: the potential benefit of using the spatial context through the AE, the input tile size and the cloud
273 cover threshold. Moreover, the spatial generalisation skill of the model is studied by splitting the colocated dataset between the
274 Northern and Southern hemispheres. For each of these, the performance for the AT variant of the OR model is reported as it
275 performs significantly better than the IT variant across experiments and evaluation metrics.

276

277 2.5.1 Spatial context

278 In order to evaluate the actual effect of the spatial context with respect to the input cloud properties, the prediction skill of the
279 model trained based on the AE encodings is compared to two trivial methods (predicting the majority bin and predicting the bin
280 minimising the MAE across the training dataset) and a method relying on the flattened cloud properties of a 9x9 tile centred
281 around the observation. Both of the trivial methods result in always predicting the CBH bin of 600 m. The third method yields a
282 similar dimensionality as the AE encodings (3 channels x 9 x 9 = 243) and thus helps to show how the AE potentially leverages
283 some spatial information about the cloud scene. Across all metrics, the baseline method using the 9x9 tile input is outperformed
284 by the initial method and even by the trivial choice of the majority bin, increasing the MA-RMSE by 400 m and the MA-MAE
285 by 140 m compared to the OR predictions made with the AE. Using the trivial choice of the 600 m bin results in an increase of
286 the MA-MAE (+7.7%) and of the MA-RMSE (+4.8%) compared to the base method. The mean bias of the trivial method is
287 lowered closer to 0 m as it leads to a more substantial underestimation of the high CBHs and overestimation of the low CBHs. To
288 conclude the comparison with these two other baselines, the information spatially encoded by the AE over the whole tile size

289 area is useful in producing CBH retrievals of better quality compared to a baseline OR model with a reduced spatial context or a
290 trivial method predicting a singular bin.

291

292 2.5.2 Tile size

293 A prediction model is fitted to the input data using encodings produced with tailored AE models trained as detailed in the
294 previous section but with varying square input tile sizes of 16, 64 and 128. With the subsequent prediction models, the retrievals
295 made with a tile size of 128 showcase the lowest MA-MAE (0.8% and 2.7% decreases compared to tile sizes of 16 and 64
296 respectively) and MA-RMSE (around a 5% decrease compared to both other tile sizes), while no clear sensitivity arises from the
297 OC, UOC or AUOC. Examining performance for each class separately indicates reduced errors (MAE and RMSE) for higher
298 CBHs (above 1000 m) using the larger tile size of 128 and on par performance across tile sizes for lower CBHs. In the context of
299 the presented CBH retrieval, the larger spatial information provided through the input tile seems to be useful for the subsequent
300 CBH prediction task, leveraged with the help of the AE as shown previously.

301

302 2.5.3 Cloud cover

303 The colocated dataset is first filtered again with cloud cover thresholds of 10%, 20% and 30%. Each threshold respectively leads
304 to datasets of 25 042, 23 034 and 21 065 samples which are then further splitted in training, validation and testing. On the
305 validation set, while the decreases in MA-MAE (4.5%) and MA-RMSE (10%) with the 10% compared to the 30% cloud cover
306 threshold are indicating a potential benefit of lowering the threshold, investigating the MAE and class-wise MAEs sheds a
307 different picture: the benefit seems to marginally concern the higher CBH classes while hindering performances on low CBHs
308 which overall explains the trend in RMSE notably. Considering the confusion matrices generated for each cloud cover threshold
309 additionally shows that a lower cloud cover threshold results in a slightly increasing distribution shift of the predicted CBH
310 classes towards higher CBHs, displaying a prediction cluster around 1000m. Overall, the benefit of additional available samples
311 when lowering the cloud cover threshold does not seem to directly lead to convincing improved performance. The main axis of
312 improvement here is probably lying in the widening of the collocation process to ensure broader spatial and temporal coverage of
313 the training dataset.

314

315 2.5.4 Spatial generalisation

316 Furthermore, in a similar way as for investigating the spatial generalisation ability of the AE, we split our colocated dataset
317 between the Northern and Southern hemispheres. This way, we ensure a minimal amount of samples in each spatial split (17 615
318 and 3 450 for the Northern and Southern hemispheres respectively) even though the spatial distribution patterns of the retrievals
319 greatly differ. As a result, the lower amount of samples in the Southern hemisphere leads to some overfitting with metrics
320 systematically worsening when testing on the Northern hemisphere. However, the Northern hemisphere training displays fair
321 generalisation skill with equal or improved metrics when testing on the Southern hemisphere, for example an 8% decrease in
322 MA-RMSE, 1% decrease in OC and stable MA-MAE, UOC and AUOC. The class-wise performances for the two splits reveal
323 the overall generalisation difficulty for higher CBHs (above 600 m) when training on the Southern hemisphere, as the labels
324 relative to these classes are mostly present in the Northern hemisphere (Figure A.3). The ability of the model to generalise from
325 the Northern hemisphere labels reassures the overall skill of the model once trained on all the labels available.

326

327 In the following section, we present the results of the developed method alongside comparisons to previous retrieval approaches.
328 In particular, we compare our retrieval to a method assuming an adiabatic cloud model (adapted from Goren et al. (2018), cf.
329 appendix E for implementation) and to the method from Noh et al. (2017). The former relies on the CTH retrieved from
330 CALIPSO's Cloud Aerosol Lidar with Orthogonal Polarization (CALIOP; Hunt et al., 2009) and CloudSat (Stephens et al.,
331 2008), but CWP and CTT retrievals from MODIS MYD06. However, in our own comparison study we used all necessary
332 variables, including the CTH, from MODIS MYD06. The latter method relies on piecewise linear relationships between MODIS
333 CWP and the geometric thickness of the uppermost layer from CALIPSO/CloudSat stratified by MODIS CTH. The application
334 of the method presented in Noh et al. (2017) is however done with CTH retrievals from the Suomi-National Polar-Orbiting
335 Partnership (SNPP) VIIRS. The comparison to our method presented here is done by using the
336 MODIS/CALIPSO/CloudSat-derived parameters from Noh et al. (2017), but using the MODIS derived CTH to produce the final
337 CBH estimate. In both cases, since these methods can be applied pixel-wise when a MODIS retrieval is available, we computed
338 the retrieved CBH values and averaged them over the cloud scene.

339 3 Results, evaluation, and comparison to previous retrieval approaches

340

341 3.1 Cloud base height retrieval, evaluation and comparison to previous retrievals

342

343 In this section, we present the results of the retrieval, evaluate it using the ground-based observations, and investigate how our
344 method fares by comparing it to a method assuming an adiabatic cloud model (adapted from Goren et al. (2018), cf. appendix E
345 for implementation) and to the method from Noh et al. (2017). The analysis is performed for the co-located scenes where
346 ground-based observations are available. To be able to compare the relevant metrics for the different methods we proceed to a
347 binning of the data following the WMO standard presented in section 2.1. In [Table 2](#) we report several metrics including the
348 MAE, the mean error (bias), the RMSE and the standard deviation of the absolute error. The latter helps us characterise the
349 spread and uncertainty in the overall predictions with respect to the surface observations. We additionally report the adapted
350 version of the AUOC mentioned in section 2.5. Furthermore, we do not report quantities such as the correlation coefficient or the
351 regression line on the 2-dimensional histograms of [Figure 5](#) and [Figure 6](#), as the stratified and categorical aspects of the data
352 would make reporting these not clearly informative. We refer to the overall conceived method including the AE (cf. section 2.4)
353 and the OR prediction model in the AT variant (cf. section 2.5), listed in [Table 2](#) as ORABase.

354 We first note that the OR method with an immediate-threshold setup fails at predicting the cloud scene base height with similar
355 skill compared to the other retrieval products, producing large errors (double-fold in comparison to the all-threshold setup). On
356 the other hand, ORABase performs well with satisfying error measures and uncertainty in the predictions on par with the other
357 retrievals. Compared to the method from Noh et al. (2017), our method succeeds in decreasing on average the error, displaying a
358 reduction of 100 m for the MAE. The method also effectively diminishes the uncertainty in the CBH retrievals, bringing down
359 the absolute error standard deviation 200 m lower. Our method thus provides accurate retrievals with comparatively low general
360 uncertainty levels. Even though on average the predictions exhibit a slight positive bias, we find that the CBH values above
361 2000 m are systematically underestimated ([Fig. 5](#)). In consideration of the low representation of such observations in the dataset,
362 due to data filtering and surface observations being less reliable for higher clouds, the method still struggles to properly quantify
363 the cloud scene base height of these samples. These samples also make up for most of the measurement uncertainty in the labels
364 considering that ceilometers face challenges for retrieving cloud signals higher up in the boundary layer. Focusing on lower
365 cloud scene base height retrievals, the predictions demonstrate even lower errors: the MAE is lowered to 379 m while the
366 absolute error standard deviation is narrowed down to 328 m. Achieved accuracy levels and uncertainty measures attest to a
367 certain trustworthiness of the cloud scene base height estimates, in particular in the context of product requirements for example
368 the ones outlined by the Joint Polar Satellite System (JPSS; Goldberg et al. (2013); 2 km accuracy threshold). However, the cloud
369 scene base height retrieval method presented here does not aim at constituting a product on its own as it is not operational with
370 the processing of daily new data available from the MODIS instrument, but rather at providing robust estimates of CBH for
371 lower level clouds. Therefore, it is expected and reasonable that the accuracies and uncertainties presented here are below such
372 thresholds. However, the available method code (Lenhardt et al., 2024) easily allows the processing of new data for users, in
373 addition to the available dataset for the year 2016.

374 We performed further sensitivity studies on our retrieval method trying to improve the quality of the predictions. An attempt to
375 balance the dataset by oversampling the higher CBH values (cloud base retrievals falling into the 2500 m bin), however, did not
376 yield better results overall but also posed a higher risk of overfitting to these specific samples. Furthermore, any spatial
377 information about the location of the satellite retrieval was not included as to prevent possible overfitting to the latitude and
378 longitude coordinates of the observations present in the training data. Since the observations are sparsely distributed especially in
379 the southern hemisphere (cf. figures from appendix A), the goal is to avoid any kind of induced spatial bias and sensitivity in the
380 model's predictions. Accordingly we can then ensure proper generalisation skill to new spatial areas, but not only based on
381 known retrieval distributions at similar locations. As a consequence, the choice was made to evaluate the potential generalisation
382 skill of the prediction model by establishing a geographic distribution of the mean predicted cloud scene base height for a whole
383 year's worth of MODIS overpasses. This is discussed in more detail in section 4. On the other hand, the temporal aspect of the
384 model's generalisation skill was intrinsically ensured by building a test set temporally distinct from the training set, including
385 co-located samples only from the last months of 2016.

386

387

388

Method	MAE (m)	Bias (m)	RMSE (m)	Absolute error standard deviation (m)	AUOC
Goren et al. (2018)	457	- 262	689	515	0.92
Noh et al. (2017)	578	- 35	860	638	0.92
OR (IT) + AE	991	+ 595	1296	836	0.93
ORABase	447	+ 58	614	420	0.89
ORABase training	456	+ 80	620	420	0.89

389

390 **Table 2: Performance on the test set of different CBH retrieval methods. OR models are either built with the**
391 **immediate-threshold (IT) or all-threshold (AT) variant. The method on which the rest of the study is based has been**
392 **highlighted in bold and its corresponding performance on the training set is added in the last row.**

393

394 3.2 Comparison to spaceborne radar-lidar retrievals of the CBH

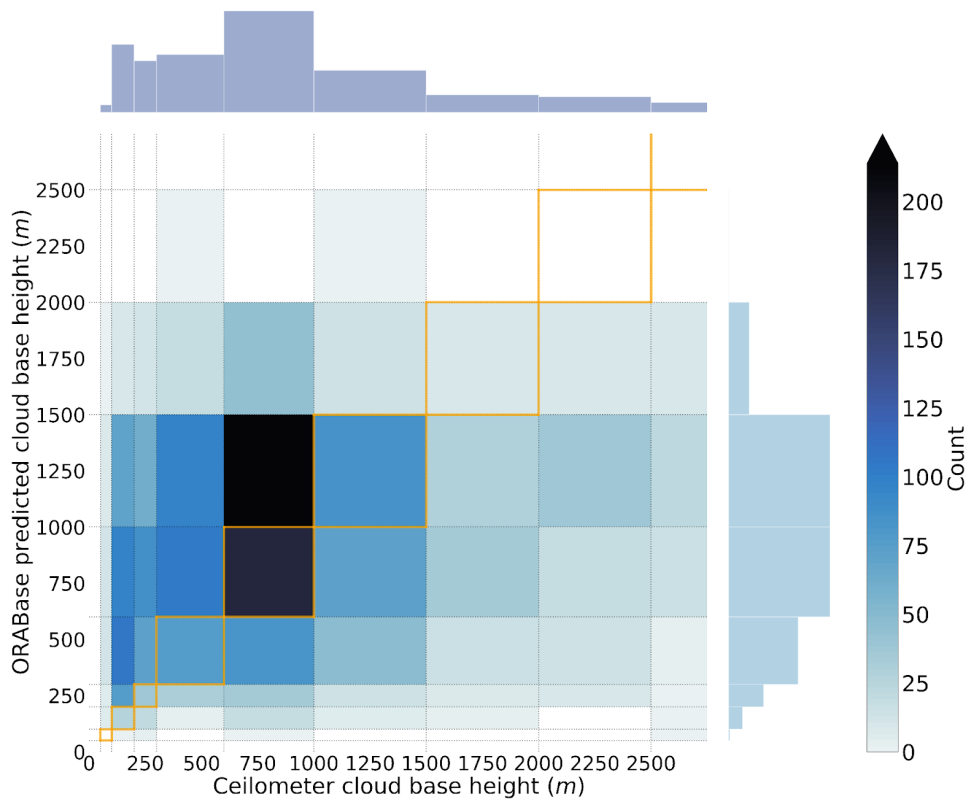
395

396 The combined datasets which are part of CUMULO (Zantedeschi et al., 2019), in particular the radar and lidar retrievals,
397 facilitate the joint evaluation of our method with both ceilometer surface observations and active satellite retrievals. Specifically
398 we leverage the 2B-CLDCLASS-LIDAR product (Sassen et al., 2008) which is derived from the combination of CloudSat’s
399 Cloud Profiling Radar (CPR; Stephens et al., 2008) and CALIPSO’s Cloud-Aerosol Lidar with Orthogonal Polarisation
400 (CALIOP; Hunt et al., 2009). The base height of the lowest cloud layer retrieved by the instruments in each scene is considered
401 the scene CBH and then averaged over the available pixels along the track, preserving the same spatial extent as the associated
402 cloud properties from the MODIS instrument. For the co-located samples of the year 2008, we thus jointly retrieve the obtained
403 CBH from the 2B-CLDCLASS-LIDAR product, only considering cases where a surface observation was in the vicinity of the
404 satellite track (inside a disc with a ~60 km radius around the surface observation, cf. section 2.3). For the samples fulfilling these
405 conditions, we then compare how the different retrievals fare. In [Figure 6](#), the joint histograms for the surface observations, the
406 2B-CLDCLASS-LIDAR retrieval and the method’s corresponding predictions are documented, representing a total of around
407 800 samples.

408 Investigating the joint histogram between the surface observations and the 2B-CLDCLASS-LIDAR retrievals ([Fig. 6a](#)) allows to
409 identify shortcomings of the active satellite retrievals in particular close to the surface (Tanelli et al., 2008; Marchand et al.,
410 2008). Indeed, the CBHs closer to the surface are not well captured by the 2B-CLDCLASS-LIDAR retrievals as partially
411 expected, due to thick clouds attenuating the lidar signal, and due to ground clutter and lack of sensitivity to small droplets near
412 cloud base for the radar signal. A similar explanation can eventually be articulated as a whole for the co-located retrievals,
413 considering that the mean bias between the two retrievals is greater than + 600 m. Concurrently, it is fruitful to compare the
414 2B-CLDCLASS-LIDAR retrievals with the predictions from the developed method ([Fig. 6b](#)). As seen previously, ORABase
415 struggles at higher CBHs, but agrees here reasonably well with the active satellite retrievals, especially for retrievals between
416 500 m and 1500 m. Focusing on retrievals under 1.5 km, the prediction model achieves similar performance as presented in [Table](#)
417 [2](#) with a MAE of 488 m and a RMSE of 576 m, even though the subset here is much smaller.

418 Furthermore, we created a more extensive dataset using only 2B-CLDCLASS-LIDAR retrievals and the cloud scene predictions
419 with the aim of obtaining a more complete view of the relationship between these two retrievals. To this extent, we collated
420 around 160 000 samples of aligned cloud scene base height predictions and the 2B-CLDCLASS-LIDAR retrievals over the year
421 2016. For this dataset, the performance metrics exhibit similar values as on the previously presented subset, displaying even
422 lower values for the MAE and the absolute error standard deviation (around a 50 m decrease for both). Similarly to the previous
423 co-located subset, limiting the evaluation to lower cloud base retrievals yields performance metrics close to a 450 m MAE and a
424 270 m absolute error standard deviation, both of these being mainly impacted by agreeing retrievals in the 500 m to 1500 m
425 range.

Joint histogram - Surface observations and model predictions



426

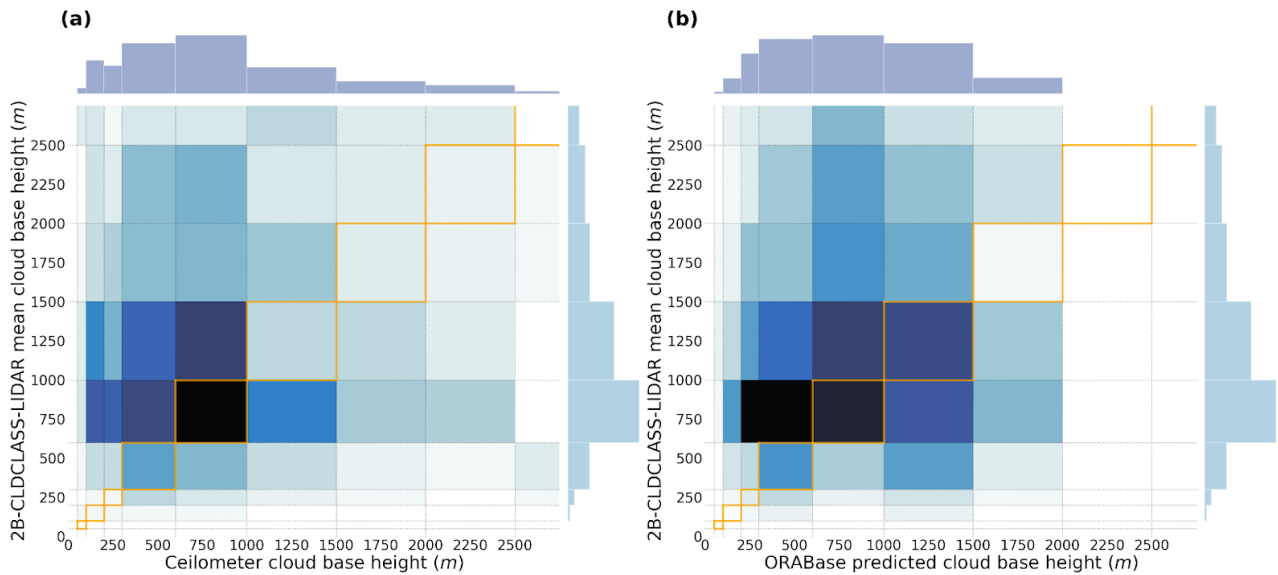
427

428 **Figure 5: Joint histogram over the test set of the surface observations and the predicted cloud scene base height from**
429 **ORABase with the ordinal regression all-threshold model. The 1:1 boxes are highlighted in orange in the figure.**

430 4 Global distribution

431

432 To further evaluate the method, we also apply the prediction model on global MODIS data for the whole year of 2016. The
433 sampling process yields approximately 700 000 CBH retrievals for the corresponding cloud properties tiles. We then spatially
434 aggregate the predictions to a regular grid of 5° and compute the annual mean per grid cell along the annual median absolute
435 deviation (MAD). The MAD constitutes a useful metric to quantify the variability while removing the effects of outliers. For
436 more robust evaluation and statistics, only ocean grid cells with more than 100 CBH retrievals over the year are displayed thus
437 impacting mostly coastal and polar regions where filtering for ocean-only scenes or the original amount of satellite retrievals
438 leads to a higher rate of displaying removal. The spatial distribution of the mean cloud base (Fig. 7, top) is similar to the outlined
439 global distributions from other studies using different instruments and methods (Böhm et al., 2019; Lu et al., 2021; Mülmenstädt
440 et al., 2018). The illustrated global quantities were established using MODIS overpasses which happen at a practically constant
441 local time (13:30 h, early afternoon for AQUA). The MAD pattern exhibits similar characteristics (Fig. 7, bottom), even though
442 variability slightly increases in the vicinity of land masses. These interpretations still remain valid when looking at relative
443 deviations. Typical features are lower cloud bases towards polar regions and the mid-latitudes, and higher ones in the tropical
444 regions. One can further observe regions like the pacific coast of South America or the Namibian coast which display lower
445 cloud bases concurrently with lower variability (also highlighted in Lu et al. (2021)). It is however impossible to follow up the
446 study for nighttime retrievals, as some MODIS cloud properties are not retrieved then.



447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

Figure 6: Joint histogram of (a) surface observations and 2B-CLDCLASS-LIDAR retrievals, and (b) ORABase predictions and 2B-CLDCLASS-LIDAR retrievals, for the co-located cloud scenes during the year 2008. The 1:1 boxes are highlighted in the figure in orange.

5 Conclusion

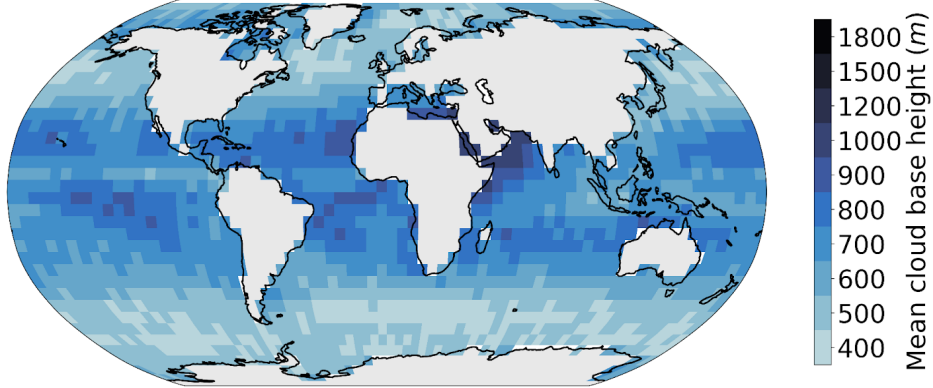
We have presented here a novel method named ORABase which retrieves the cloud scene base height over marine areas from MODIS cloud properties, specifically CTH, COT and CWP. This method can produce robust CBH estimates for cloud scenes in particular for lower cloud bases (MAE of 379 m and absolute error standard deviation of 328 m for up to 2 km cloud bases), based on the assumption of a homogeneous cloud base across the considered cloud field. The statistical model was built on surface observations of cloud bases with ceilometers (section 2.1), and then evaluated in comparison to other methods using passive satellite instruments (section 3.1) and active satellite retrievals (section 3.2). Analysis of the yearly averaged CBH (section 4) helped to further make sense of the predicted cloud bases and variability. The global dataset for the year 2016 is available from Zenodo (Lenhardt et al., 2024).

Using the spatially-resolved information of cloud fields of CTH, COT and CWP through the described CNN-AE results in more accurate CBH retrievals compared to the active retrievals of the 2B-CLDCLASS-LIDAR product, producing better performance metrics compared to the other products and methods considered in this study. The combination of a CNN based AE to reduce the dimensionality of the spatial patterns of cloud properties followed by a simple OR model leads to a better CBH retrieval compared to previous presented methods. The OR modelisation helps bridging the gap between regression and classification, facilitating the use of the binned cloud base observations provided by the surface observation dataset. Overall, ORABase achieves low error in the retrievals, around 400 m, and concurrently a narrow absolute error distribution, more precisely around 400 m absolute error standard deviation. Both of these performance metrics are additionally reduced when focusing on cloud bases lower than 2 km. Application to data over land areas has not been processed yet but would certainly require adding surface observations from land during the training process (e.g. Böhm et al., 2019; Lu et al., 2021; Mülmenstädt et al., 2018). Application of the presented retrieval method to other instruments could also be considered. Incorporating TERRA MODIS data would help constrain the annual mean estimates presented in Figure 7 by partially removing the potential bias of the single daily overpass arising from using only AQUA data presented in this study. The aspect enabling potential application of the retrieval method to different instruments outside of the two MODIS sensors would be the standardisation process for the input cloud properties before the use of the AE which is done based on means and standard deviations computed from AQUA-only granules. Carefully investigating the characteristics of the distribution of the cloud properties from another instrument to ensure proper scaling when using the trained AE would be then necessary. Further tests could be additionally done using coarser resolution for the input cloud properties.

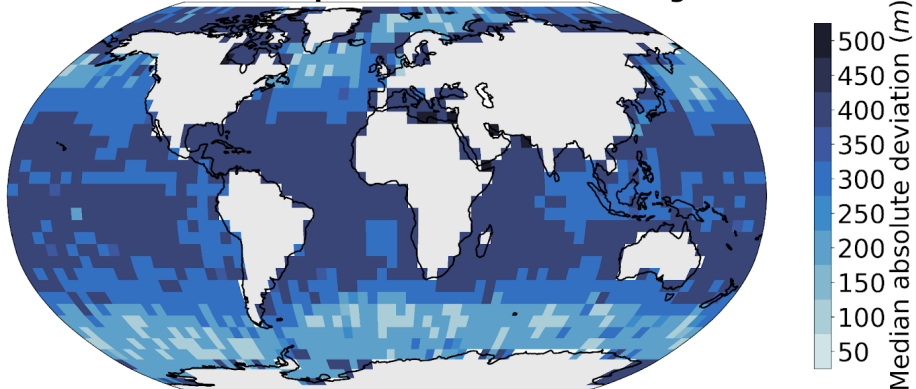
Furthermore, classical semi-supervised pipelines like the one presented here, characterised by a small labelled dataset and a vast unlabelled dataset, necessitate a kind of co-location or matching process which often proves to be cumbersome and generates only a limited amount of labels. However, future avenues of research could consider directly modelling unmatched datasets, as in

483 e.g. Lun Chau et al. (2021) with multiresolution atmospheric data, by making use of other quantities present in the observations
484 as mediating variables to model the link between observed and unobserved variables.
485 In essence, the main benefit of producing better cloud base estimates is to gain accuracy in the overall retrieval of cloud
486 geometry, impacting in particular radiation estimates (Kato et al., 2011) like the surface downwelling longwave radiation
487 (Mülmenstädt et al., 2018). ORAbase can thus prove to be useful by helping to produce CBH with enhanced confidence at a
488 global scale.
489
490

Mean predicted cloud base height - Year 2016

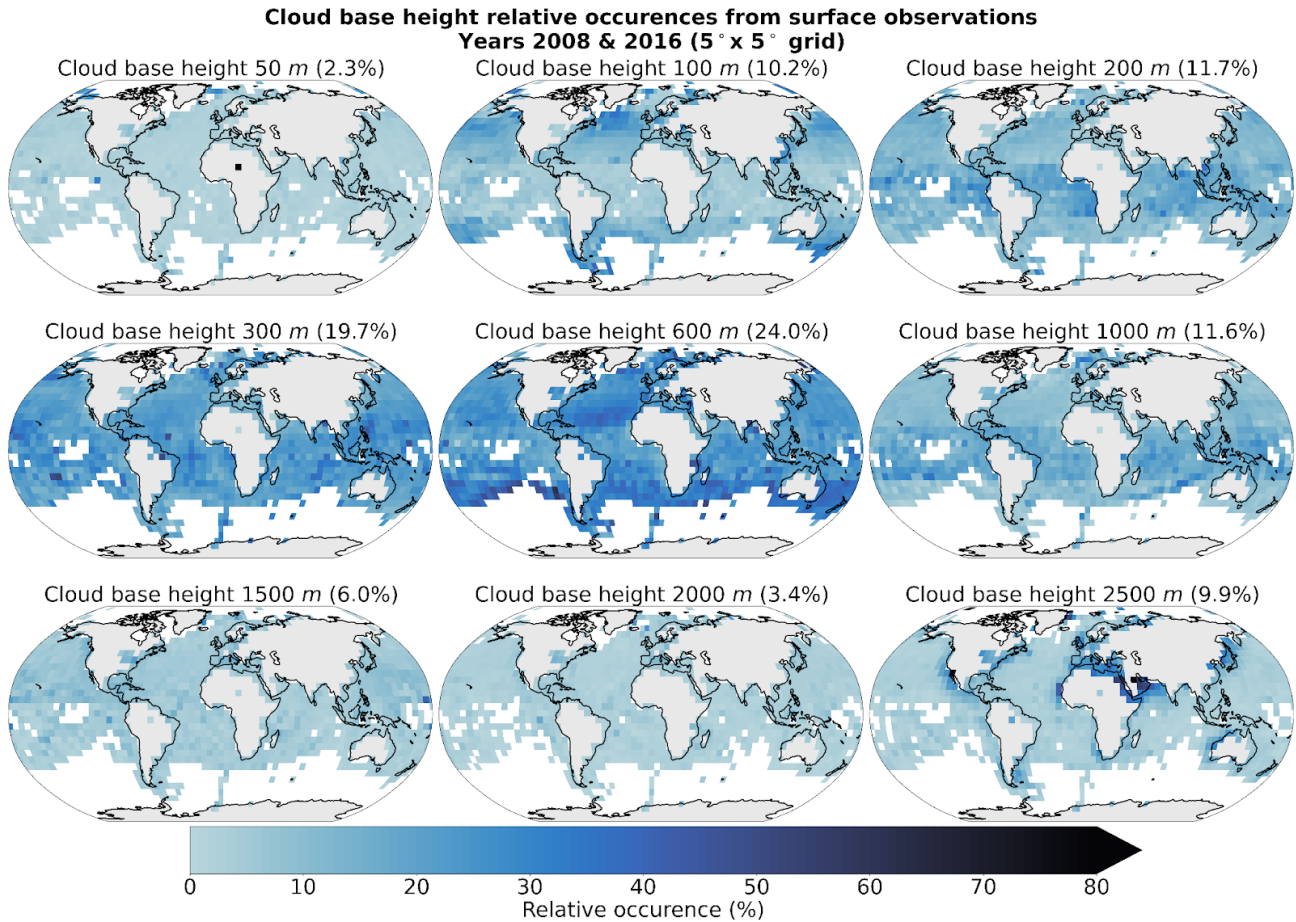


Median absolute deviation of predicted cloud base height - Year 2016



491
492 **Figure 7: Spatial distribution of (top) mean and (bottom) median absolute deviation of predicted cloud base height for the**
493 **MODIS data of the year 2016 aggregated on a 5 ° grid.**

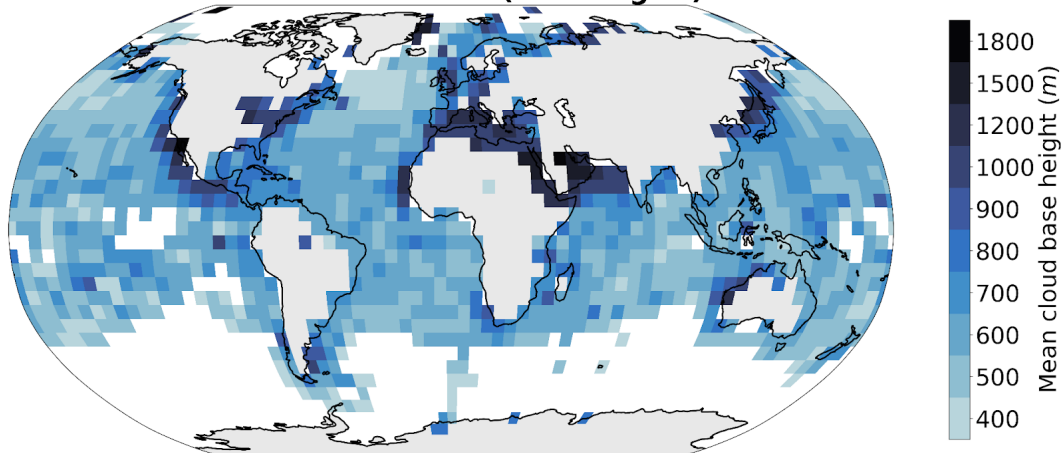
496 Appendix A: Cloud base height retrievals distribution



497

498 **Figure A.1: Spatial distribution of cloud base height retrievals (Met Office, 2006) for the years 2008 and 2016 on a 5°**
499 **grid. Overall percentage of each label in the total observations is indicated in brackets. Only grid cells with more than 10**
500 **retrievals are displayed.**

Mean cloud base height from surface observations
Years 2008 & 2016 (5° x 5° grid)

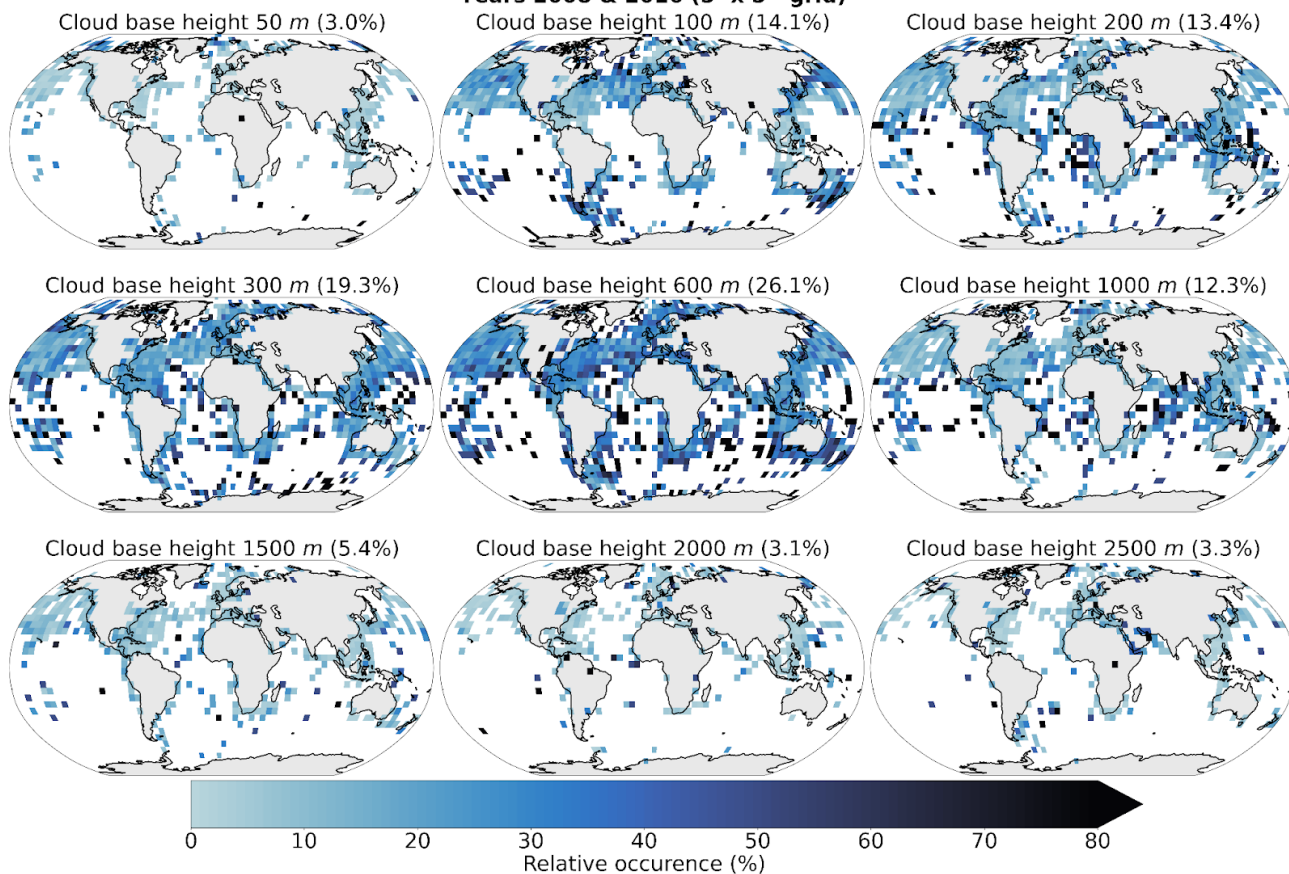


501

502 **Figure A.2: Mean cloud base height from retrievals (Met Office, 2006) for the years 2008 and 2016 on a 5° grid. Only**
503 **grid cells with more than 50 retrievals are displayed.**

504

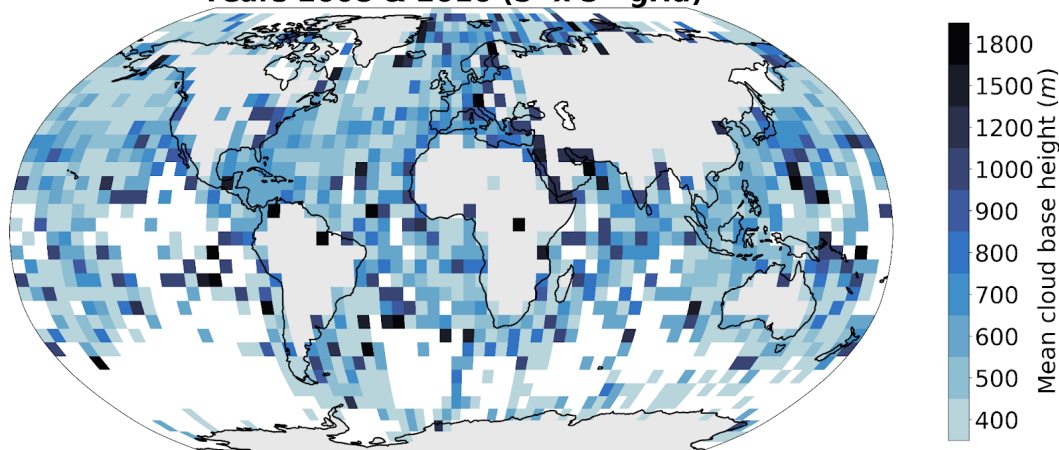
**Cloud base height relative occurrences for the co-located input dataset
Years 2008 & 2016 (5° x 5° grid)**



505
506
507
508
509

Figure A.3: Spatial distribution of the co-located cloud base height retrievals (Met Office, 2006) and the satellite cloud properties used for training the prediction model for the years 2008 and 2016 on a 5° grid. Overall percentage of each label in the total dataset is indicated in brackets.

**Mean cloud base height for the co-located input dataset
Years 2008 & 2016 (5° x 5° grid)**



510
511
512
513

Figure A.4: Mean cloud base height from the co-located retrievals (Met Office, 2006) and the satellite cloud properties used for training the prediction model for the years 2008 and 2016 on a 5° grid.

514 Appendix B: Spatio-temporal correlation study

515

516 We create five different datasets to evaluate how the chosen AE architecture is capable of generalising to new data while trying
 517 to remove some possible autocorrelation biases which might inflate the performance scores. We also use this study to analyse
 518 how the AE model behaves when trained with our input data. We define two splits for space and time in order to build the
 519 training and testing datasets, namely the South-western (SW) quadrant and the period from March to October, respectively. The
 520 granules used to build the datasets span across the whole year of 2016. The *random* data split is the basis for the training of the
 521 model and consists of tiles sampled in the aforementioned quadrant and time period. These tiles are then split randomly between
 522 training, validation and testing datasets. This split represents the common way of splitting data when building a ML model. In
 523 contrast, we build 3 other datasets which vary through their respective spatial and time spans. The *spatial* split is built
 524 considering tiles spanning across a distinct time period, here between November and February, regardless of their spatial
 525 location. The *temporal* split is built considering tiles located anywhere but in the South-western quadrant regardless of the time
 526 at which the retrieval occurred. Finally the *spatio-temporal* split combines the previous two conditions in order to build a dataset
 527 in which the tiles come from an independent location and time as the ones used for training. Additionally, we create a global data
 528 split using data from a different year, here 2008, without any spatial restriction for the tiles. Furthermore, only a limited number
 529 of tiles was extracted from each granule while only granules from non-consecutive days were used in order to limit possible
 530 correlation between the extracted scenes.

Data split	Time period	Spatial extent	<i>n</i>
Random	03-10.2016	SW quadrant	Train: 14 691 Validation: 4 198 Test: 2 099
Spatial	03-10.2016	Global except SW quadrant	107 736
Temporal	01-02 and 11-12.2016	SW quadrant	12 420
Spatio-temporal	01-02 and 11-12.2016	Global except SW quadrant	30 659
Global	12.2008	Global	7 111

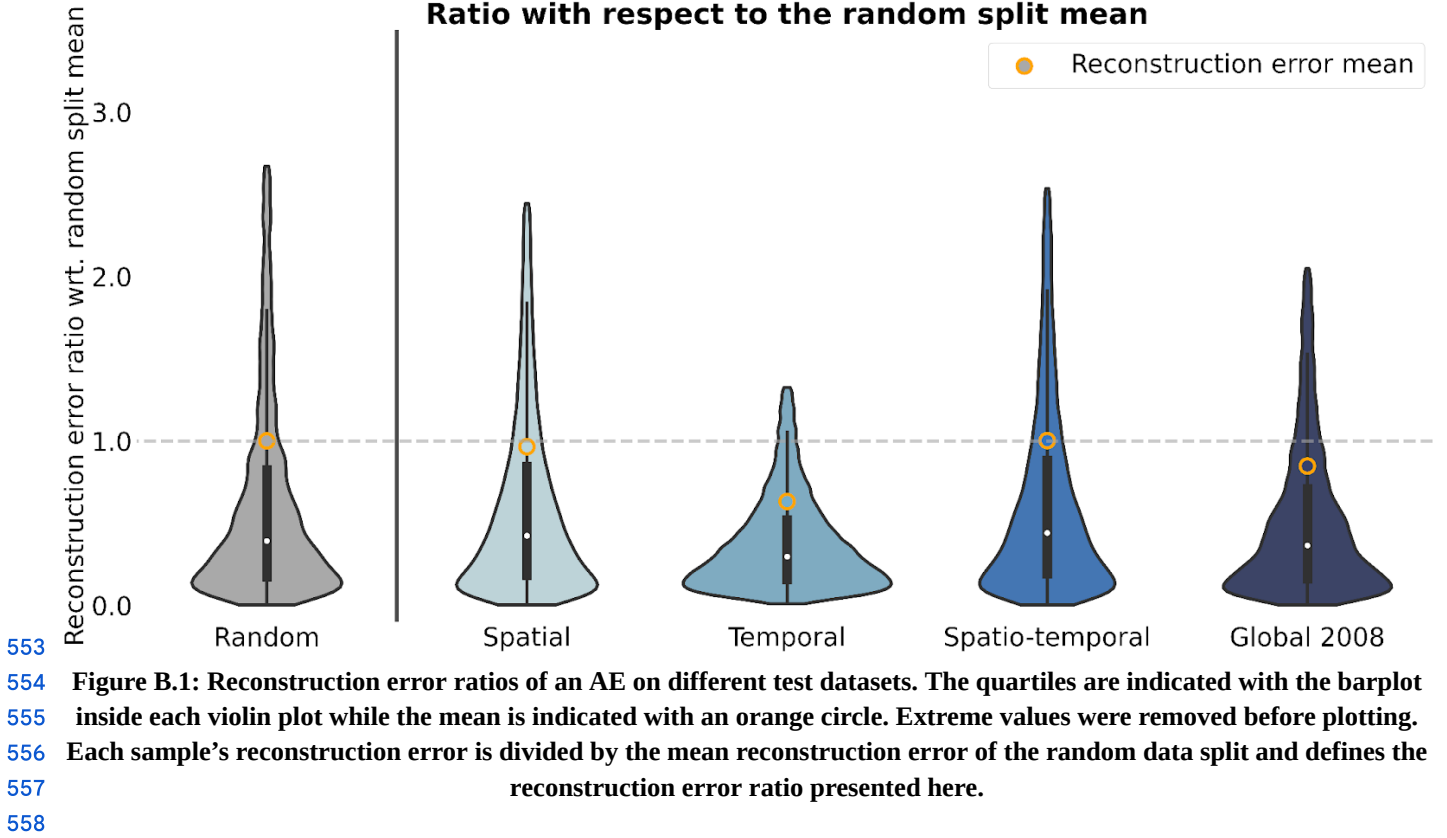
531 **Table B.1 : Name, time period, spatial extent and number of samples for each of the five described data splits.**

532

533 We then train an AE model using the training data from the first data split (*random*). Each test data split is then used to evaluate
 534 the trained model through the reconstruction errors divided by the reconstruction error mean of the *random* split (noted as
 535 reconstruction error ratio; [Fig. B.1](#)). Spatial distribution of the mean reconstruction errors is shown in [Figure B.2](#). We detail in
 536 [Table B.2](#) the average channel reconstruction error for each of the splits.

537 We first notice that the reconstruction power of the model is consistent regardless of the test split considered with mean
 538 reconstruction error ratios ranging from 0.63 to 1.0, dividing the split's reconstruction error by the random data split mean
 539 reconstruction error. Ratios around 1 or below indicate that the model's performance is not inflated when considering a random
 540 data split, highlighting that the model did not only learn from possible spatial and/or temporal correlations between samples
 541 present in the training set. The distribution of the error is also very similar throughout the test splits with most of the samples
 542 located below an error ratio of 0.5. However, one of the main aspects regarding the performance of the model across test splits is
 543 the presence of a heavy tail in the distribution showcasing that for some samples the reconstruction error can be greater than 3
 544 times the mean error. Looking at the spatial patterns of the reconstruction error, we note that overall the error comes from the
 545 COT and CWP predictions, the average reconstruction errors across test sets being 0.15, 0.32 and 0.25 for CTH, COT and CWP
 546 respectively ([Table B.2](#)). For the CTH, the error is concentrated in the zones with frequent convection around the equator and
 547 could be explained by local convection cells exhibiting a larger spread in CTH values. Another source of error could be that
 548 higher CTH values are also less represented in the training data. On the contrary, the error for COT and CWP is prevailing in
 549 high-latitude regions. Overall, the performance skill of the AE model seems to hold through the different test data splits. One
 550 could argue that the training dataset already retains enough variability in the data which could explain why the model still
 551 performs well regardless of the test set split. However, this consistent skill also shows that the performance reported in appendix
 552 C on the test set can be trusted to hold for other datasets and supports the data generation process to train the AE (cf. section 2.4).

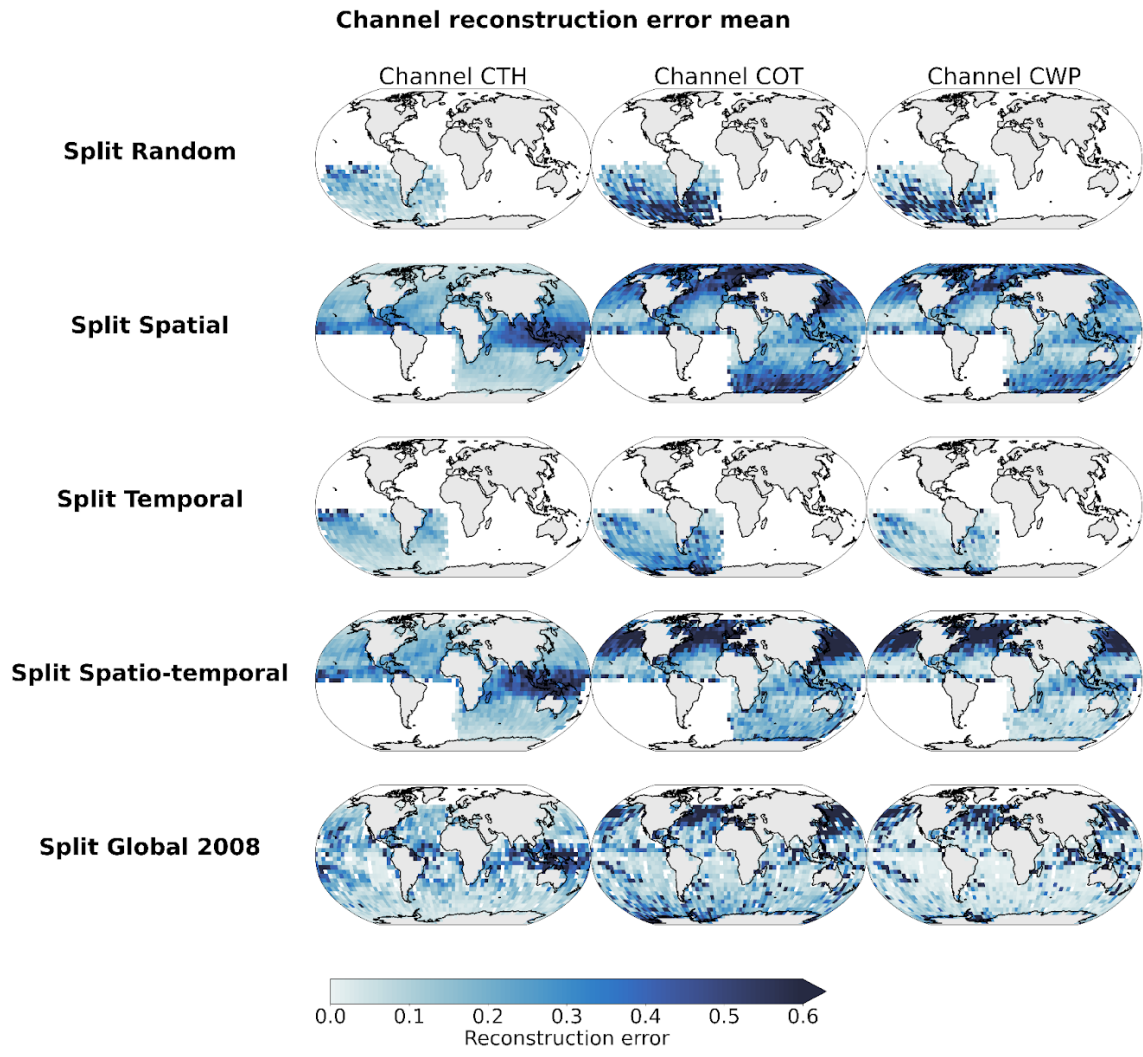
**Autoencoder reconstruction error distribution on different test sets
Ratio with respect to the random split mean**



Data split	Channel			Average
	CTH	COT	CWP	
Random	0.117	0.369	0.333	0.273
Spatial	0.171	0.344	0.276	0.263
Temporal	0.114	0.253	0.150	0.172
Spatio-temporal	0.202	0.332	0.286	0.274
Global	0.154	0.318	0.221	0.231
Average	0.152	0.323	0.253	0.243

Table B.2 : Average channel reconstruction error for each of the five described data splits.

559
560



561
562

Figure B.2: Distribution of mean channel reconstruction errors aggregated on a 5 ° grid.

Layer	Hyperparameters	Output shape
Input		(None, 3, 128, 128)
Encoder		
Conv2d	(kernel = 3, stride = 2)	(None, 3, 64, 64)
ConvBlock x 5	Conv2d (kernel = 3, stride = 1) LeakyReLU Conv2d (kernel = 3, stride = 1) LeakyReLU Conv2d (kernel = 3, stride = 1) BatchNorm2d LeakyReLU MaxPool2d (kernel = 2, stride = 2)	(None, 256, 2, 2)
Flatten + Linear		(None, 256)
Decoder		
Linear + Unflatten		(None, 256, 2, 2)
ConvTranspose2d	(kernel = 2, stride = 2)	(None, 256, 4, 4)
ConvTransposeBlock x 5	Conv2d (kernel = 3, stride = 1) LeakyReLU Conv2d (kernel = 3, stride = 1) LeakyReLU Conv2d (kernel = 3, stride = 1) BatchNorm2d LeakyReLU ConvTranspose2d (kernel = 2, stride = 2)	(None, 3, 128, 128)

Table C.1 : Autoencoder model specifications.

Hyperparameter	Value
Batch size	64
Epochs	80
Optimizer	Stochastic Gradient Descent (SGD), momentum = 0.9, learning rate = 0.0001
Metric	MSE
Early stopping	patience = 20

Table C.2 : Autoencoder model training specifications.

569 **Appendix D: Ordinal regression**

570

571 We define our labels y which can take values in $K = 9$ classes from $\{50 \text{ m}, 100 \text{ m}, \dots, 2500 \text{ m}\}$. We introduce $K - 1$
 572 thresholds α_y to define the separation of our K classes which actually correspond here to the classes too. For each labelled
 573 sample (s, y) the output of our model is $z = z(s)$. The correct interval for this sample is then (α_{y-1}, α_y) . During the fitting
 574 process, the goal is to find the set of parameters of our model z and the corresponding thresholds α which minimises a certain
 575 cost function. We consider a generic nonnegative penalisation function $f(\cdot)$ (eg. hinge loss, squared error loss, Huber loss).
 576 There are then different ways to represent threshold violations and thus to penalise the predictor. While immediate-threshold
 577 setup only considers the thresholds of the correct interval, all-threshold setup takes into account all the threshold violations. In
 578 the case of an immediate-threshold setup the loss function would look like:

579
$$\mathcal{L}(z, y) = f(z - \alpha_{y-1}) + f(\alpha_y - z) \quad (\text{D.1})$$

580 Here we can see that the loss is not aware of how many thresholds are actually violated. In the case of an all-threshold setup the
 581 loss function is a sum of violations across all thresholds:

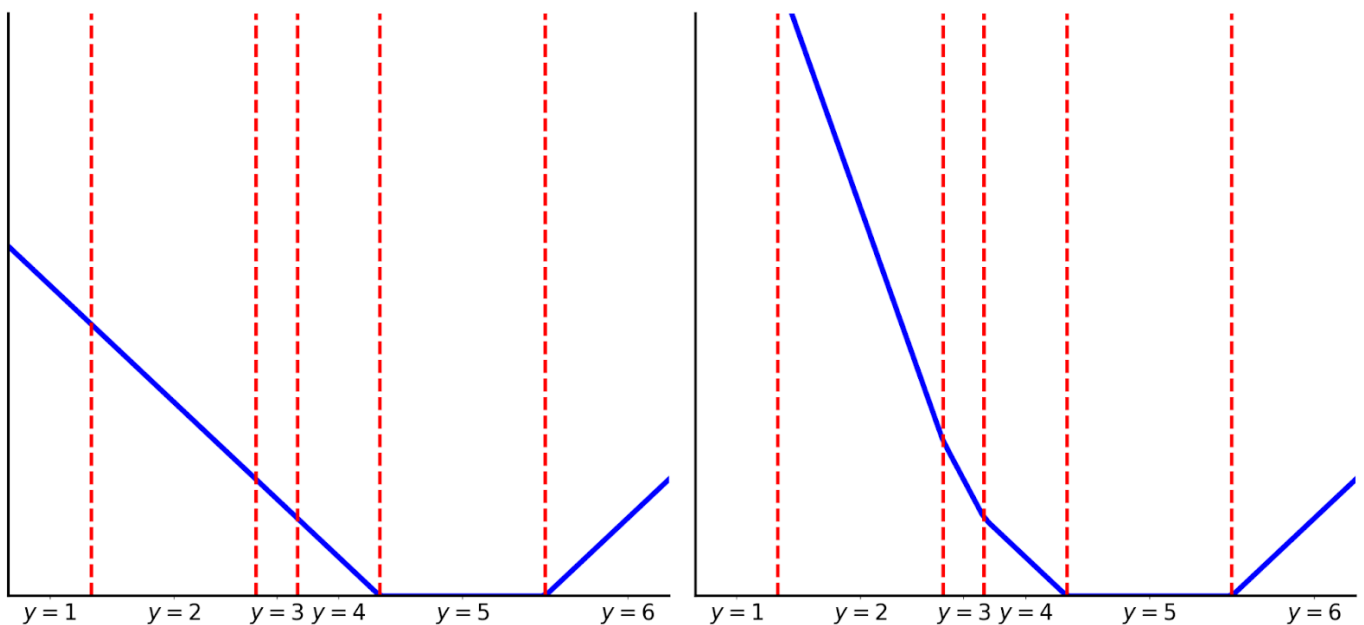
582
$$\mathcal{L}(z, y) = \sum_{i=1}^{K-1} f(t(i, y)(\alpha_i - z)) \quad (\text{D.2})$$

583 where $t(i, y) = -1$ if $i < y$ or $+1$ if $i \geq y$. Thus predictions are encouraged to violate the least amount of thresholds.

584 We give in [Figure D.1](#) an example of what the loss function would look like in the case of $K = 6$ labels and using a hinge
 585 penalisation.

586

587



588

589 **Figure D.1: Threshold-based setups loss function representation for a hinge penalisation, $K=6$ labels and target label $y=5$.**
 590 **(left) Immediate-threshold and (right) All-threshold setup loss function. (figure adapted from Rennie et al. (2005))**

591 **Appendix E: Cloud base height retrieval method assuming adiabatic cloud**

592

593 Algorithm adapted from Goren et al. (2018). We use the retrieved CTH, CTT, CTP and CWP from MODIS MYD06 (Platnick et
594 al., 2017).

595

Algorithm: Cloud base height retrieval

Data: CTH, CTT, CTP, LWP, look-up tables

Result: CBH

if CTT < 263.13 **then**

return NaN

T \leftarrow CTT - 273.13

LWP obs \leftarrow LWP

LWP adi \leftarrow 0.

$\delta z \leftarrow$ 0.

Set corresponding cloud top indexes for temperature T_{ind} and pressure p_{ind} look-up tables.

Read-in the water mixing ratio w at the corresponding indexes.

if w out of look-up table **then**

return NaN

while LWP adi < LWP obs **then**

$\rho_{tmp} \leftarrow$ density look-up table with T_{ind} and p_{ind}

$\delta_{tmp} \leftarrow$ layer depth look-up table with T_{ind} and p_{ind}

$\delta z \leftarrow \delta z + \delta_{tmp}$

$w_{tmp} \leftarrow$ mixing ratio look-up table with T_{ind} and p_{ind}

 LWP adi \leftarrow LWP adi + $(w_{tmp} - w) \times \delta z_{tmp} \times \rho_{tmp}$

 Adjust temperature T given the saturated lapse rate using look-up table with T_{ind} and p_{ind}

 Update indexes T_{ind} and p_{ind}

return CTH - δz

596

597 **Table E.1: Pseudo code for cloud base height retrieval algorithm assuming adiabatic cloud, adapted from Goren et al.**
598 **(2018).**

599 Code availability

600
601 The code used for the method and producing the plots is available on Zenodo (Lenhardt et al., 2024).

602 Data availability

603
604 The global dataset of the cloud base height predictions for the year 2016 is available on Zenodo (Lenhardt et al., 2024). The
605 dataset is available as a csv file with corresponding coordinates, MODIS granule, time of retrieval and predicted cloud base
606 height or in a netCDF file as daily aggregates on a regular grid with a resolution of 1° or 5°. The meteorological observations
607 from the UK MetOffice (Met Office, 2006) are available through the CEDA archive at
608 <https://catalogue.ceda.ac.uk/uuid/77910bcec71c820d4c92f40d3ed3f249>. The files from the CUMULO dataset (Zantedeschi et
609 al., 2019) are available at <https://www.dropbox.com/sh/i3s9q2v2jyyk2it/AACxXnXfMF5wuIqLXqH4NJOra?dl=0>.

610 Author contribution

611
612 JL, JQ and DS designed the study. JL wrote the code. JL conducted the analysis and JL, JQ, DS interpreted the results. JL
613 prepared the manuscript, JQ and DS reviewed the manuscript and provided comments.

614 Competing interests

615
616 The authors declare that they have no conflict of interest.

617 Acknowledgements

618
619 This work was supported by the European Union's Horizon 2020 research and innovation programme under Marie
620 Skłodowska-Curie grant agreement No. 860100 (iMIRACLI). We thank the Leipzig University Scientific Computing cluster for
621 computing and data hosting. We further thank Tom Goren for providing access to code snippets from Goren et al. (2018) and
622 thank Olivia Linke for helping review the manuscript. We acknowledge the contributors of the CUMULO dataset (Zantedeschi et
623 al., 2019) for providing access to the data files hosted at
624 <https://www.dropbox.com/sh/i3s9q2v2jyyk2it/AACxXnXfMF5wuIqLXqH4NJOra?dl=0>. Additionally, we acknowledge the
625 MODIS L2 Cloud product data set from the Level-1 and Atmosphere Archive and Distribution System (LAADS) Distributed
626 Active Archive Center (DAAC), located in the Goddard Space Flight Center in Greenbelt, Maryland
627 (https://ladsweb.modaps.eosdis.nasa.gov/archive/allData/61/MYD06_L2/). We would like to thank two anonymous reviewers for
628 their constructive and detailed comments.

629 References

630

631 Ackerman, S. A., and Frey, R.: MODIS Atmosphere L2 Cloud Mask Product (35_L2), NASA MODIS Adaptive Processing
632 System, Goddard Space Flight Center, http://doi.org/10.5067/MODIS/MOD35_L2.061,
633 http://doi.org/10.5067/MODIS/MYD35_L2.061, 2017.

634

635 Baccianella, S., Esuli, A. and Sebastiani, F.: Evaluation Measures for Ordinal Regression, Ninth International Conference on
636 Intelligent Systems Design and Applications, Pisa, Italy, 283-287, <https://doi.org/10.1109/ISDA.2009.230>, 2009.

637

638 Baldi, P.: Autoencoders, Unsupervised Learning, and Deep Architectures, in: Proceedings of the International Conference on
639 Machine Learning (ICML), Workshop on Unsupervised and Transfer Learning, Proceedings of Machine Learning Research,
640 Volume 27, 37-49, <https://proceedings.mlr.press/v27/baldi12a.html>, 2012.

641

642 Baum, B.A., Menzel, W. P., Frey, R. A., Tobin, D. C., Holz, R. E., Ackerman, S. A., Heidinger, A. K., and Yang, P.: MODIS
643 Cloud-Top Property Refinements for Collection 6, Journal of Applied Meteorology and Climatology, 51, 6, 1145-1163,
644 <https://doi.org/10.1175/JAMC-D-11-0203.1>, 2012.

645

646 Böhm, C., Sourdeval, O., Mülmenstädt, J., Quaas, J., and Crewell, S.: Cloud base height retrieval from multi-angle satellite data,
647 Atmos. Meas. Tech., 12, 1841-1860, <https://doi.org/10.5194/amt-12-1841-2019>, 2019.

648

649 Boucher, O., Randall, D., Artaxo, P., Bretherton, C., Feingold, G., Forster, P., Kerminen, V.-M., Kondo, Y., Liao, H., Lohmann,
650 U., Rasch, P., Satheesh, S. K., Sherwood, S., Stevens, B. and Zhang, X. Y.: Clouds and aerosols, Climate Change 2013: The
651 Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on
652 Climate Change, 571-657, <https://doi.org/10.1017/CBO9781107415324.016>, 2013.

653

654 Cardoso, J. S. and Sousa, R.: Measuring the performance of ordinal classification, International Journal of Pattern Recognition
655 and Artificial Intelligence, Volume 25, 8, 1173-1195, <https://doi.org/10.1142/S0218001411009093>, 2011.

656

657 Forster, P., T. Storelvmo, K. Armour, W. Collins, J.-L. Dufresne, D. Frame, D.J. Lunt, T. Mauritsen, M.D. Palmer, M. Watanabe,
658 M. Wild, and H. Zhang: The Earth's Energy Budget, Climate Feedbacks, and Climate Sensitivity, in Climate Change 2021: The
659 Physical Science Basis. Contribution of Working Group I to the Sixth Assessment Report of the Intergovernmental Panel on
660 Climate Change [Masson-Delmotte, V., P. Zhai, A. Pirani, S.L. Connors, C. Péan, S. Berger, N. Caud, Y. Chen, L. Goldfarb, M.I.
661 Gomis, M. Huang, K. Leitzell, E. Lonnoy, J.B.R. Matthews, T.K. Maycock, T. Waterfield, O. Yelekçi, R. Yu, and B. Zhou (eds.)].
662 Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, pp. 923–1054,
663 <http://doi.org/10.1017/9781009157896.009>, 2021.

664

665 Goldberg, M. D., Kilcoyne, H., Cikanek, H., and Mehta, A.: Joint Polar Satellite System: The United States next generation
666 civilian polar-orbiting environmental satellite system, J. Geophys. Res. Atmos., 118, 13,463–13,475,
667 <https://doi.org/10.1002/2013JD020389>, 2013.

668

669 Goren, T., Rosenfeld, D., Sourdeval, O., and Quaas, J.: Satellite Observations of Precipitating Marine Stratocumulus Show
670 Greater Cloud Fraction for Decoupled Clouds in Comparison to Coupled Clouds, Geophys. Res. Lett., 45, 5126–5134,
671 <https://doi.org/10.1029/2018GL078122>, 2018.

672

673 Grosvenor, D. P., Sourdeval, O., Zuidema, P., Ackerman, A., Alexandrov, M. D., Bennartz, R., Boers, R., Cairns, B., Chiu, J. C.,
674 Christensen, M., Deneke, H., Diamond, M., Feingold, G., Fridlind, A., Hünerbein, A., Knist, C., Kollias, P., Marshak, A.,
675 McCoy, D., Merk, D., Painemal, D., Rausch, J., Rosenfeld, D., Russchenberg, H., Seifert, P., Sinclair, K., Stier, P., van
676 Diedenhoven, B., Wendisch, M., Werner, F., Wood, R., Zhang, Z. and Quaas, J.: Remote sensing of droplet number concentration
677 in warm clouds: A review of the current state of knowledge and perspectives, Reviews of Geophysics, 56, 409–453,
678 <https://doi.org/10.1029/2017RG00059>, 2018.

679

680 Gutiérrez, P. A., Pérez-Ortiz, M., Sánchez-Monedero, J., Fernández-Navarro, F. and Hervás-Martínez, C.: Ordinal Regression
681 Methods: Survey and Experimental Study, *IEEE Transactions on Knowledge and Data Engineering*, 28, 1, 127-146,
682 <https://doi.org/10.1109/TKDE.2015.2457911>, 2016.

683

684 Hinton, G.E., and Salakhutdinov, R.R.: Reducing the dimensionality of data with neural networks, *Science*, 313, 5786, 504-507,
685 <https://doi.org/10.1126/science.1127647>, 2006.

686

687 Hunt, W. H., Winker, D. M., Vaughan, M. A., Powell, K. A., Lucker, P. L., and Weimer, C.: CALIPSO Lidar Description and
688 Performance Assessment. *J. Atmos. Oceanic Technol.*, 26, 1214–1228, <https://doi.org/10.1175/2009JTECHA1223.1>, 2009.

689

690 Ioffe, S., and Szegedy, C.: Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift, in:
691 *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, *Proceedings of Machine Learning Research*,
692 Volume 37, 448-456, <http://proceedings.mlr.press/v37/ioffe15.html>, 2015.

693

694 Kato, S., Rose, F. G., Sun-Mack, S., Miller, W. F., Chen, Y., Rutan, D. A., Stephens, G. L., Loeb, N. G., Minnis, P., Wielicki, B.
695 A., Winker, D. M., Charlock, T. P., Stackhouse, P. W. J., Xu, K.-M., and Collins, W. D.: Improvements of top-of-atmosphere and
696 surface irradiance computations with CALIPSO-, CloudSat-, and MODIS-derived cloud and aerosol properties, *J. Geophys.*
697 *Res.-Atmos.*, 116, D19209, <https://doi.org/10.1029/2011JD016050>, 2011.

698

699 Kattenborn, T., Schiefer, F., Frey, J., Feilhauer, H., Mahecha, M. D., and Dormann, C. F.: Spatially autocorrelated training and
700 validation samples inflate performance assessment of convolutional neural networks, *ISPRS Open Journal of Photogrammetry*
701 *and Remote Sensing*, 5, 2667-3932, <https://doi.org/10.1016/j.ophoto.2022.100018>, 2022.

702

703 Kramer, M.A.: Nonlinear principal component analysis using autoassociative neural networks, *AIChe J.*, Volume 37, 233-243,
704 <https://doi.org/10.1002/aic.690370209>, 1991.

705

706 Krizhevsky, A., Sutskever, I., and Hinton, G.: ImageNet Classification with Deep Convolutional Neural Networks, in:
707 *Proceedings of Advances in Neural Information Processing Systems 25*, Annual Conference on Neural Information Processing
708 Systems (NeurIPS), 1097-1105,
709 https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf, 2012.

710

711 Lázaro, M, and Figueiras-Vidal, A. R.: Neural network for ordinal classification of imbalanced data by minimizing a Bayesian
712 cost, *Pattern Recognition*, Volume 137, <https://doi.org/10.1016/j.patcog.2023.109303>, 2023.

713

714 LeCun, Y., Jackel, L. D., Boser, B., Denker, J. S., Graf, H. P., Guyon, I., Henderson, D., Howard, R. E., and Hubbard, W.:
715 Handwritten digit recognition: Applications of neural network chips and automatic learning, *IEEE Communications Magazine*,
716 Volume 27, Issue 11, 41-46, <https://doi.org/10.1109/35.41400>, 1989.

717

718 LeCun, Y., and Bengio, Y.: Convolutional networks for images, speech, and time series, *The handbook of brain theory and neural*
719 *networks*, 3361, 10, 1995.

720

721 LeCun, Y., Kavukcuoglu, K., and Faret, C.: Convolutional networks and applications in vision, in: *Proceedings of 2010 IEEE*
722 *International Symposium on Circuits and Systems*, 253-256, <https://doi.org/10.1109/ISCAS.2010.5537907>, 2010.

723

724 Lenhardt, J., Quaas, J., and Sejdinovic, D.: Code and data for: Marine cloud base height retrieval from MODIS cloud properties
725 using machine learning, Zenodo, <https://doi.org/10.5281/zenodo.10517687>, 2024.

726

727 Lu, X., Mao, F., Rosenfeld, D., Zhu, Y., Pan, Z., and Gong, W.: Satellite retrieval of cloud base height and geometric thickness of
728 low-level cloud based on CALIPSO, *Atmos. Chem. Phys.*, Volume 21, Issue 15, 11979-12003,
729 <https://doi.org/10.5194/acp-21-11979-2021>, 2021.

730

731 Lun Chau, S., Bouabid, S., and Sejdinovic, D.: Deconditional Downscaling with Gaussian Processes, in: Proceedings of
732 Advances in Neural Information Processing Systems 34, Annual Conference on Neural Information Processing Systems
733 (NeurIPS), <https://doi.org/10.48550/arXiv.2105.12909>, 2021.
734

735 Maas, A. L., Hannun, A. Y. and Ng, A. Y.: Rectifier Nonlinearities Improve Neural Network Acoustic Models, in: Proceedings of
736 the 30th International Conference on Machine Learning (ICML), Atlanta, Georgia, USA, Journal of Machine Learning Research
737 (JMLR), Volume 28, 3, 2013.
738

739 Marchand, R., Mace, G. G., Ackerman, T., and Stephens, G.: Hydrometeor detection using Cloudsat – An earth-orbiting 94-GHz
740 cloud radar, *J. Atmos. Ocean. Technol.*, Volume 25, 519–533, <https://doi.org/10.1175/2007JTECHA1006.1>, 2008.
741

742 Met Office: MIDAS: Global Marine Meteorological Observations Data, NCAS British Atmospheric Data Centre,
743 <https://catalogue.ceda.ac.uk/uuid/77910bcec71c820d4c92f40d3ed3f249>, 2006.
744

745 Mülmenstädt, J., Sourdeval, O., Henderson, D. S., L'Ecuyer, T. S., Unglaub, C., Jungandreas, L., Böhm, C., Russell, L. M., and
746 Quaas, J.: Using CALIOP to estimate cloud-field base height and its uncertainty: the Cloud Base Altitude Spatial Extrapolator
747 (CBASE) algorithm and dataset, *Earth System Science Data*, Volume 10, Issue 4, 2279–2293,
748 <https://doi.org/10.5194/essd-10-2279-2018>, 2018.
749

750 Nair, V., and Hinton, G. E.: Rectified linear units improve restricted boltzmann machines, in: Proceedings of the 27th
751 International Conference on International Conference on Machine Learning (ICML'10), Haifa, Israel, 2010, 807–814,
752 <https://www.cs.toronto.edu/%7Efritz/absps/reluICML.pdf>, 2010.
753

754 Niu, Z., Zhou, M., Wang, L., Gao, X., and Hua, G.: Ordinal Regression with Multiple Output CNN for Age Estimation, IEEE
755 Conference on Computer Vision and Pattern Recognition (CVPR), 4920-4928, <https://doi.org/10.1109/CVPR.2016.532>, 2016.
756

757 Noh, Y., Forsythe, J. M., Miller, S. D., Seaman, C. J., Li, Y., Heidinger, A. K., Lindsey, D. T., Rogers, M. A., and Partain, P. T.:
758 Cloud-Base Height Estimation from VIIRS. Part II: A Statistical Algorithm Based on A-Train Satellite Data, *Journal of*
759 *Atmospheric and Oceanic Technology*, Volume 34, Issue 3, 585-598, <https://doi.org/10.1175/JTECH-D-16-0110.1>, 2017.
760

761 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison,
762 A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J. and Chintala, S.:
763 PyTorch: An Imperative Style, High-Performance Deep Learning Library, in *Advances in Neural Information Processing*
764 *Systems* 32 (NeurIPS), 8024–8035,
765 <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>, 2019.
766

767 Pedregosa, F.: Feature extraction and supervised learning on fMRI: from practice to theory, Université Pierre et Marie Curie,
768 Paris VI, <https://theses.hal.science/tel-01100921>, 2015.
769

770 Pedregosa, F., Bach, F., and Gramfort, A.: On the Consistency of Ordinal Regression Methods, *Journal of Machine Learning*
771 *Research (JMLR)*, Volume 18, 55, 1-35, <http://jmlr.org/papers/v18/15-495.html>, 2017.
772

773 Platnick, S., Ackerman, S. A., King, M. D., Meyer, K., Menzel, W. P., Holz, R. E., Baum, B. A., and Yang, P.: MODIS
774 atmosphere L2 cloud product (06_L2), NASA MODIS Adaptive Processing System, Goddard Space Flight Center,
775 http://doi.org/10.5067/MODIS/MYD06_L2.061, 2017.
776

777 Platnick, S., King, M.D., Ackerman, S.A., Menzel, W.P., Baum, B.A., Riedi, J.C., and Frey, R.A.: The MODIS cloud products:
778 algorithms and examples from Terra, in: *IEEE Transactions on Geoscience and Remote Sensing*, Volume 41, Number 2, 459-473,
779 <http://doi.org/10.1109/TGRS.2002.808301>, 2003.
780

781 Pu, Y., Gan, Z., Henaio, R., Yuan, X., Li, C., Stevens, A., and Carin, L.: Variational Autoencoder for Deep Learning of Images,
782 Labels and Captions, in: Proceedings of Advances in Neural Information Processing Systems 29, Annual Conference on Neural
783 Information Processing Systems (NeurIPS), 2352-2360, <https://doi.org/10.48550/arXiv.1609.08976>, 2016.
784

785 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., and Prabhat: Deep learning and process
786 understanding for data-driven Earth system science, *Nature*, 566, 195-204, <https://doi.org/10.1038/s41586-019-0912-1>, 2019.
787

788 Rennie, J.D., and Srebro, N.: Loss Functions for Preference Levels : Regression with Discrete Ordered Labels, in: Proceedings of
789 the IJCAI multidisciplinary workshop on advances in preference handling, Volume 1, 180–186, AAAI Press, Menlo Park, CA,
790 2005.
791

792 Ronneberger, O., Fischer, P., and Brox, T.: U-Net: Convolutional Networks for Biomedical Image Segmentation, in: Navab, N.,
793 Hornegger, J., Wells, W., Frangi, A. (eds) *Medical Image Computing and Computer-Assisted Intervention (MICCAI 2015)*,
794 *Lecture Notes in Computer Science*, Volume 9351, Springer, Cham., https://doi.org/10.1007/978-3-319-24574-4_28, 2015.
795

796 Sassen, K., Wang, Z., and Liu, D.: Global distribution of cirrus clouds from CloudSat/Cloud-Aerosol Lidar and Infrared
797 Pathfinder Satellite Observations (CALIPSO) measurements, *J. Geophys. Res.*, Volume 113, D00A12,
798 <https://doi.org/10.1029/2008JD009972>, 2008.
799

800 Shi, X., Cao, W., and Raschka, S.: Deep Neural Networks for Rank-Consistent Ordinal Regression Based On Conditional
801 Probabilities, *Pattern Analysis and Applications*, Volume 26, 941–955, <https://doi.org/10.1007/s10044-023-01181-9>, 2023.
802

803 Silva, W., Pinto, J. R., and Cardoso, J. S.: A Uniform Performance Index for Ordinal Classification with Imbalanced Classes,
804 2018 International Joint Conference on Neural Networks (IJCNN), Rio de Janeiro, Brazil, 1-8,
805 <https://doi.org/10.1109/IJCNN.2018.8489327>, 2018.
806

807 Simonyan, K., and Zisserman, A.: Very Deep Convolutional Networks for Large-Scale Image Recognition, 3rd International
808 Conference on Learning Representations (ICLR), Computational and Biological Learning Society, 1-14,
809 <https://arxiv.org/abs/1409.1556>, 2015.
810

811 Stephens, G. L., Vane, D. G., Tanelli, S., Im, E., Durden, S., Rokey, M., Reinke, D., Partain, P., Mace, G. G., Austin, R.,
812 L'Ecuyer, T., Haynes, J., Lebsock, M., Suzuki, K., Waliser, D., Wu, D., Kay, J., Gettelman, A., Wang, Z., and Marchand, R.:
813 CloudSat mission: Performance and early science after the first year of operation, *J. Geophys. Res.*, Volume 113, D00A18,
814 <http://doi.org/10.1029/2008JD009982>, 2008.
815

816 Tanelli, S., Durden, S. L., Im, E., Pak, K. S., Reinke, D. G., Partain, P., Haynes, J. M., and Marchand, R. T.: CloudSat's Cloud
817 Profiling Radar After Two Years in Orbit: Performance, Calibration, and Processing, *IEEE Trans. Geosci. Remote Sens.*, Volume
818 46, 3560–3573, <https://doi.org/10.1109/TGRS.2008.2002030>, 2008.
819

820 TorchVision maintainers and contributors: TorchVision: PyTorch's Computer Vision library, GitHub repository,
821 <https://github.com/pytorch/vision>, 2016.
822

823 Trenberth, K. E., Fasullo, J. T., and Kiehl, J.: Earth's global energy budget, *Bulletin of the American Meteorological Society*,
824 Volume 90, 311–324, <http://doi.org/10.1175/2008BAMS2634.1>, 2009.
825

826 Watson-Parris, D., Rao, Y., Olivié, D., Seland, Ø., Nowack, P., Camps-Valls, G., Stier, P., Bouabid, S., Dewey, M., Fons, E.,
827 Gonzalez, J., Harder, P., Jeggle, K., Lenhardt, J., Manshausen, P., Novitasari, M., Ricard, L., and Roesch, C.: ClimateBench v1.0:
828 A benchmark for data-driven climate projections, *Journal of Advances in Modeling Earth Systems*, Volume 14, Issue 10,
829 <https://doi.org/10.1029/2021MS002954>, 2022.
830

831 Winship, C., and Mare, R. D.: Regression Models with Ordinal Variables, *American Sociological Review*, Volume 49, Number 4,
832 512–525, <https://doi.org/10.2307/2095465>, 1984.

833

834 WMO: Manual on Codes (WMO-No. 306), Volume I.1, Part A, Alphanumeric codes, Code table 1600,
835 <https://library.wmo.int/idurl/4/35713>, 2019.

836

837 Zantedeschi, V., Falasca, F., Douglas, A., Strange, R., Kusner, M. J., and Watson-Parris, D.: Cumulo: A Dataset for Learning
838 Cloud Classes, Tackling Climate Change with Machine Learning Workshop, 33rd Conference on Neural Information Processing
839 Systems (NeurIPS 2019), Vancouver, Canada, <https://doi.org/10.48550/arXiv.1911.04227>, 2019.

840

841 Zeiler, M. D. , Krishnan, D., Taylor, G. W., and Fergus, R. : Deconvolutional networks, in: Proceedings of the 2010 IEEE
842 Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), San Francisco, CA, USA, 2528-2535,
843 <https://doi.org/10.1109/CVPR.2010.5539957>, 2010.