

---

# RESIDUAL U-NET WITH ATTENTION FOR DETECTING CLOUDS IN SATELLITE IMAGERY

---

A PREPRINT

**P Shokri**

SkyWatch Space Applications Inc.  
3 - 8 Queen Street North  
Kitchener, Ontario, Canada  
N2H 2G8  
pshokri@skywatch.com

**AL De Souza**

SkyWatch Space Applications Inc.  
3 - 8 Queen Street North  
Kitchener, Ontario, Canada  
N2H 2G8  
adesouza@skywatch.com

November 29, 2022

## ABSTRACT

Semantic segmentation of clouds in Earth observation imagery is an important task in a variety of remote sensing contexts: from the application of atmospheric corrections, to being able to accurately omit cloud pixels when extracting information about ground features. Here we introduce a deep learning approach based on the popular U-Net architecture. The core of the architecture is a U-Net with residual units that ease the training of the network. An attention mechanism is also incorporated to enable the model to more effectively learn and distinguish between cloud and non-cloud features. We also explore two complementary loss functions, Binary Cross Entropy and Jaccard, in order to overcome data imbalances common to this application. Our model is trained on a uniquely curated high resolution dataset spanning a wide variety of scene contexts, lighting conditions, and seasonality. Our experiments demonstrate that this model is an extremely accurate and robust model for the semantic segmentation of clouds in satellite imagery, and the model achieves state-of-the-art performance over many other models (including others based on CNN architectures) on common benchmark datasets, even without having been exposed to those datasets prior to testing.

**Keywords** remote sensing · cloud detection · Landsat-8 · Sentinel-2 · deep learning · U-Net · attention

## 1 Introduction

Approximately 70% of the Earth’s surface is obscured by cloud at any given time when observed from space [King et al., 2013]. This high degree of cloud cover is problematic for a variety of Earth observation applications that depend on passive remote sensing in the visible parts of the electromagnetic spectrum: image compositing [Roy et al., 2010]; the application of atmospheric corrections [Vermote et al., 2002]; calculation of vegetation indices [Huete et al., 2002]; land cover classification [Zhang et al., 2002]; and especially change detection [Zhu and Woodcock, 2014]. Accurate identification and screening of clouds in satellite imagery is therefore an essential step toward producing high-quality geoinformatic products.

A variety of methods exist for the purpose of detecting clouds in remote sensing data. Many of these solutions rely on rule-based or threshold-based methods which use the reflectance difference between cloud and non-cloud pixels on multiple spectral bands. Zhu and Woodcock [2012] proposed the Function of Mask (FMask) for detecting clouds, utilizing seven specific bands in Landsat data. FMask detects cloud pixels based on a threshold function and the physical characteristics of clouds such as whiteness, flatness, and temperature. Though several improvements have been made to the algorithm over time [Zhu et al., 2015, Qiu et al., 2019], in practice, threshold based methods such as FMask are not reliable in isolation and often require additional information, such as surface temperature. Li et al. [2017] proposed automatic multi-feature combined cloud detection (MFC). This method uses low-level features, such

as textural and geometrical features across Red, Green, Blue (RGB), and near-infrared (NIR) bands. However both FMask and MFC are limited in extracting high-level semantic cloud features, and are prone to detection errors.

Recently, deep learning models have begun to outperform such methods in detecting clouds. López-Puigdollers et al. [2021] demonstrated that in an inter-class experiment on Landsat-8 (L8) data, where the training and inference sets are independent but from the same satellite source, deep learning models have considerably lower commission errors compared to FMask, while having similar accuracy scores. The authors compared several cloud segmentation models trained on L8 imagery, and tested on different combinations of L8 and Sentinel-2 (S2) data. One of the models in their benchmark is the lightweight U-Net proposed by Zhang et al. [2020]. They find that although threshold based methods can perform well in detecting thin clouds, deep learning based approaches outperform such models in distinguishing low contrast differences, such as cloud from snow pixels.

More recently Chen et al. [2021] proposed an Automatic Cloud Detection neural network (ACD net) to semantically segment clouds. ACD net integrates geospatial data and satellite imagery, together with a novel cloud boundary refinement module, to increase the accuracy of cloud detection by their network. The authors used the Gaofen-1 satellite imagery dataset ( $512 \times 512$  pixels) provided by Li et al. [2017], selecting image patches that include cloud-snow coexistence scenes as well as different types of cloud patterns. The authors found ACD net outperforms other deep learning-based models in distinguishing such scenes, but does not cope well with thin cloud layers [Li et al., 2017, Zhan et al., 2017, Xia et al., 2019, Yang et al., 2019].

Hu et al. [2021] proposed a novel variation of U-Net called CDUNet with ResNet-50 [He et al., 2016] as its encoder, designed to refine the edge of cloud segmentation regions. They performed their experiments on RGB bands of L8-SPARCS ( $256 \times 256$  pixels) and Google Earth’s cloud and cloud shadow ( $224 \times 224$  pixels) datasets. The authors proposed novel modules in the decoder to recover lost details in the encoding stage, reduce information redundancy and reconstruct spatial information using a self-attention mechanism. However, CDUNet is computationally expensive, having approximately 47.59 million parameters. Compared to the original U-Net model, CDUNet’s F1 score was an approximately 2% improvement on both datasets.

Following the Attention U-Net architecture [Oktay et al., 2018], Guo et al. [2020] introduced a similar model with a soft-attention mechanism called Cloud-AttU to also detect clouds in L8 imagery. They used a cross entropy loss function to train their model on  $384 \times 384$  pixel images with RGB and near-infrared bands. Cloud-AttU outperforms the original U-Net by 2.7% in the Jaccard index on the Landsat-Cloud dataset [Mohajerani and Saeedi, 2019], and successfully distinguishes between snow and cloud pixels. However the authors did not experiment with other inference image datasets.

In this paper we introduce a new cloud segmentation model based on an U-Net architecture with an attention mechanism [Oktay et al., 2018] as well as residual blocks in the decoder [He et al., 2016]. We compare the effectiveness of using a cross entropy loss versus Jaccard loss function [Rahman and Wang, 2016] in training the model, report the results of rigorous experiments on challenging scenarios at high spatial resolution, in mixed environmental scenes, and explore the model’s ability to generalize when tested on completely unseen datasets. These experiments demonstrate the robustness of our model to different sensors and spatial resolutions. Our model is trained and tested exclusively on RGB data. This further simplifies the architecture in comparison to other such deep-learning models, and further emphasizes the robustness of our model with respect to sensor type.

This paper is organized as follows: Section 1.1 describes our SkyWatch Cloud Segmentation dataset. Section 2 introduces the semantic segmentation model (which we refer to as Attention ResUNet) and the methodology used in this study. In Section 3 we report on the performance and results of this network on a variety of intra- and inter-class experiments. Finally, Section 4 concludes this paper with a brief summary and a discussion of potential improvements for future work.

## 1.1 SkyWatch Cloud Segmentation Dataset

The SkyWatch Cloud Segmentation (SW-CS) dataset consists of a combination of low and high spatial resolution images captured in RGB by several different sensors. In order to produce a dataset including enough variation, each of the following subgroups have been manually annotated by different sets of people. Different pre-processing techniques are then applied to each of these subgroups.

The SW-CS dataset consists of 2,576 full scene products. These products are obtained at high spatial resolution (0.5-m per pixel edge) by several independent Earth observation satellites, and come in a variety of sizes and dimensions. We use a combination of techniques to crop and rescale these products to produce subsets of imagery that reflect a variety of circumstances in which they might be used, and to allow the network to more easily process and ingest them. For 836 of the full scenes we down-sample the scene to be 2,048 pixels along their longest dimension. For 1,340 we

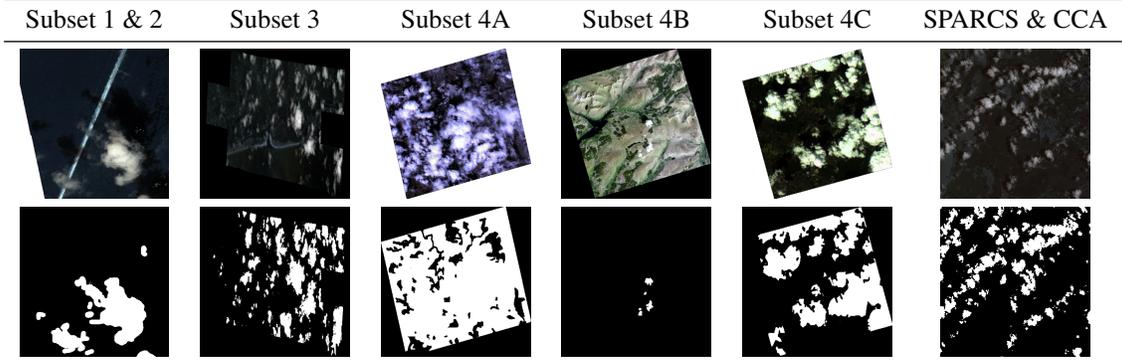


Table 1: Examples of RGB images (top row) and their corresponding cloud masks (bottom row) for each of the subsets described in the SW-CS dataset. Subset 1 & 2: Tiled 0.5-m resolution scenes. Subset 3: Rescaled 0.5-m resolution scenes. Subset 4A: 0.7-m resolution scenes. Subset 4B: 0.8-m resolution scenes. Subset 5: 2.2-m resolution scenes. SPARCS & CCA: Tiled 30-m resolution scenes.

down-sample the scenes to be 1,000 pixels in width. For the remaining 400 scenes we down-sample them to be 800 pixels wide in their shortest dimension. In the first two cases the overall aspect ratio of the imagery is preserved, while in the third case the aspect ratio is not maintained at all.

In the first two cases (the 836 + 1,340 scenes) the resized scene is then cropped into  $512 \times 512$  pixel tiles before being ingested by the model. Individual tiles with fewer than 50% valid pixels are discarded. The remaining 400 800-pixel wide scenes are simply resized to  $512 \times 512$  pixels in their entirety, and after tiling, ingested as-is regardless of the individual tiles' valid pixel count.

Working with such high resolution imagery, it is computationally expensive to ingest full sized products, which are often ten or more thousands of pixels in a single dimension. But these specific manipulations have been made to provide the model with the opportunity to gain the full context of all of the features in each scene. All of the scenes have had their cloud mask labels manually refined in order to ensure the most accurate possible cloud masks for each.

Finally, we incorporate into the SW-CS dataset a variety of lower resolution imagery as well: 134 scenes obtained at 0.7-m resolution, 54 scenes at 0.8-m resolution, and 128 scenes at 2.2-m resolution. These additional 316 scenes are all rescaled from their original size to be  $512 \times 512$  pixels, previously being between 300 and 600 pixels in each dimensions. Then additionally, the scenes in this subset have had their cloud masks manually labeled from scratch.

The SPARCS [Hughes and Hayes, 2014] and CCA [Foga et al., 2017] datasets containing L8 images and their cloud masks were also used in this study. The spatial resolution of these images is 30-meters. Each SPARCS image is split into 4 tiles (with no overlap). Each tile is then resized to be  $512 \times 512$  pixels. The result is a collection of 320 images. Each CCA image is also tiled into  $512 \times 512$  image segments so that none of the images contain invalid pixels. The result is a collection of 719 images. The original labels of these two datasets have not been refined.

### 1.1.1 Intra-class data pre-processing

The SW-CS dataset is shuffled together with the SPARCS and CCA datasets, and then split into train, test, and validation sets in a 0.8 : 0.1 : 0.1 ratio. The training set is then augmented using the following procedure:

- Rotate images in the range of  $[-20, +20]$  degrees.
- Zoom in and out in the range of  $[-20, +20]\%$ .
- Flip images in the x-direction 80% of the time.
- Flip images in the y-direction 80% of the time.
- Change the pixel intensity values in the range of  $[-20, 35]$ .
- Adjust the image contrast by scaling pixel values.  $255 \times (v/255)^\gamma$  by changing  $\gamma$  in the range  $[0.66, 1.9]$ .
- Add additive Gaussian noise per pixel, sampling from the normal distribution  $N(0, 0.06 * 255)$ . 50% of the time this noise is added to a single channel, and 50% of the time it is added equally to all three channels.

## 2 Semantic Segmentation

U-Net [Ronneberger et al., 2015] is a popular neural network model consisting of an encoder-decoder architecture, and which has gradually been improved since its original publication [Oktay et al., 2018, Piao and Liu, 2019, Cao and Zhang, 2020, Isensee and Maier-Hein, 2020]. In this paper we incorporate the suggestions provided by Oktay et al. [2018], Isensee and Maier-Hein [2020] for improving U-Net’s accuracy. We exploit an attention mechanism at the skip connections in U-Net’s decoder path. The implemented attention mechanism suppresses activations at irrelevant regions, improving the accuracy of the feature representation sent to the decoder. We further make use of ResNet-50 [He et al., 2016] pre-trained on ImageNet [Deng et al., 2009] as the decoder of our model. The skip connections in its residual blocks help address the problem of vanishing gradients (Pascanu et al., 2013; conv1\_conv, conv2\_block3\_1\_relu, and conv3\_block4\_1\_relu, in the Keras implementation for example, and as illustrated in Figure 1). The resulting encoder-decoder structure has 11.2 million trainable parameters. We therefore refer to our model as an Attention ResUNet. Figure 1 shows the architecture of our Attention ResUNet, and Figure 2 shows the structure of the attention gate used in the model.

The encoder’s weights are not frozen during training. Stochastic gradient descent is facilitated with the Adam optimizer [Kingma and Ba, 2014], and early stopping [Caruana et al., 2000] is used to prevent over-fitting during training, where the validation loss is monitored. If the loss does not decrease for 50 consecutive epochs, the training process is stopped.

The resolution of the input images to the Attention ResUNet model is  $512 \times 512$  pixels. Images consist of three wavelength bands: red, green, and blue. The model’s output is a segmentation mask, also referred to as a cloud mask, which is a single band  $512 \times 512$  pixel image. The pixel values in a cloud mask correspond to the probability  $([0, 1])$  of that pixel being all or part of a cloud. In this context we refer to the *background* versus *cloud* classes.

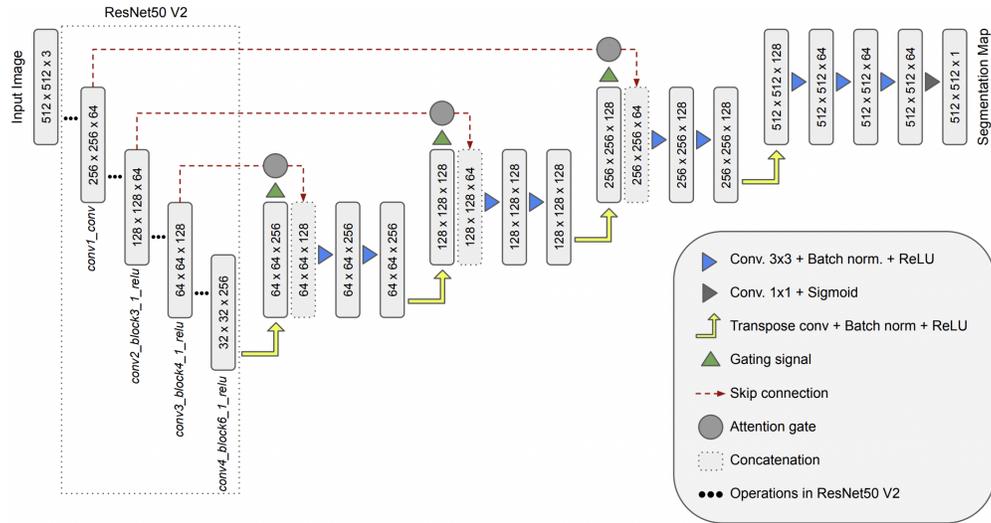


Figure 1: Block diagram of the Attention ResUNet segmentation model. The input image is filtered and down-sampled by the ResNet-50 V2 architecture in the encoder. Attention gates filter features coming from the encoder. The structure of the attention gate is shown in Figure 2 below.

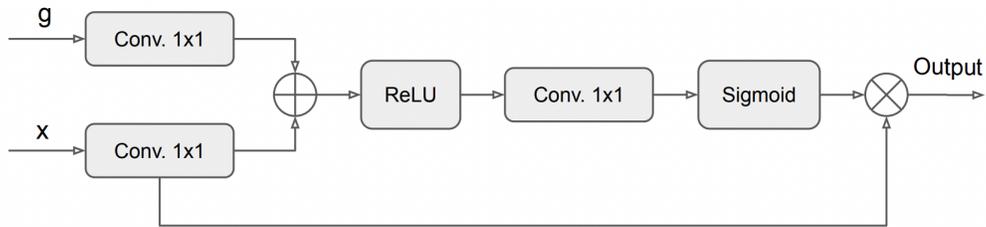


Figure 2: The structure of the attention gate in the Attention ResUNet, where  $g$  is the gating signal and  $x$  is the skip connection obtained from the encoder. The  $\oplus$  symbol represents the addition operation, and  $\otimes$ , the dot product operation.

## 2.1 Loss Functions

Janocha and Czarnecki [2017] note that cross entropy loss is often used for training neural networks without consideration for the alternatives. Indeed, cross entropy has been the fiducial loss function in applications of deep learning to cloud detection [Hu et al., 2021, Guo et al., 2020]. In the case of highly imbalanced images where the cloud pixels are few, however, this loss function can lead to simplistic solutions because many pixels can be wrongly classified as the background class to ensure a low loss value. Jaccard loss, also known as Intersection over Union (IoU) Rahman and Wang [2016], is a great candidate for dealing with highly imbalanced datasets. Therefore, we have trained three variations of our cloud detection model.

We trained one of the variations of our model using Binary Cross Entropy (BCE) loss, as in equation 1. Here  $T_x$  refers to true labels or the ground truth of pixel  $x$ , and  $P_x$  refers to the network’s prediction of pixel  $x$ .

$$L_{BCE} = \sum_x -(T_x \log(P_x) + (1 - T_x) \log(1 - P_x)) \quad (1)$$

The second variation of our model is trained using Jaccard loss, introduced in equation 2. Here  $T$  stands for the true label while  $P$  stands for the prediction. The notation  $| \cdot |$  is the 1-norm.  $T * P$  is the element-wise multiplication. The nominator of this formula is the approximate intersection based on the probabilities of  $P$  and  $T$ .

$$L_{Jaccard} = \frac{|T * P|}{|T + P - (T * P)|} \quad (2)$$

Finally, equation 3 expresses a third variation of our model, which is trained using the average of both the BCE and Jaccard losses:

$$L_{avg} = \frac{L_{BCE} + L_{Jaccard}}{2} \quad (3)$$

## 3 Results

In this section we measure the performance of our Attention ResUNet model as trained on the SW-CS dataset. First, the intra-class inference results are provided and discussed. Second, we perform an inter-class experiment (without retraining) on the various datasets mentioned in [Skakun et al., 2022]. We aim to evaluate our model against a common benchmark, and for this we follow the methodology of the Cloud Mask Inter-comparison eXercise (hereafter CMIX, Skakun et al., 2022).

### 3.1 Performance metrics

As in Skakun et al. [2022] we use overall accuracy (OA), balanced overall accuracy (BOA; Brodersen et al., 2010), producer’s accuracy (PA, or recall R), and user’s accuracy (UA, equivalent to precision P) to report the performance of models in both intra- and inter-class experiments. These metrics are calculated using the total number of true positives, false positives, and false negatives in our inference experiments. For the inter-class experiment, we additionally calculated the precision and recall per inference image and report the average performance of both, which we call  $P_{mean}$  and  $R_{mean}$ , respectively. We also report the F1 and IoU scores.  $F1_{mean}$  and  $IoU_{mean}$  being the averages of these scores on all inference images, while  $F1_{total}$  and  $IoU_{total}$  are calculated using the total number of true positives, false positives, and false negatives in the inference dataset.

### 3.2 The Effects of Loss Functions

The output of the segmentation network is a mask in which each pixel’s value represents the probability of that pixel being part of a cloud. In other words, each pixel value determines the confidence of the network in identifying cloud vs ground features. A threshold on this probability can be applied to fix each pixel as being either a cloud or non-cloud pixel. Higher threshold values result in fewer false positives (increasing precision and decreasing recall), whereas lower threshold values result in decreased precision and increased recall. Throughout our experiments below we use a threshold of 50% to ensure a balance between the model’s precision and recall.

In experimenting with both BCE and Jaccard loss we find that the cloud masks of the network when trained with Jaccard loss is strongly determined, with probabilities being always close to 0 or 1. With BCE loss the inferred cloud masks are often more ”gray”, with pixel values taking on a much more uniform distribution of values 0 and 1. Interestingly, the combination of both losses resulted in cloud masks that are not as discriminatory as using the Jaccard

loss alone, nor is there as much variance in their output as there is with BCE by itself. An example is provided in Figure 3.

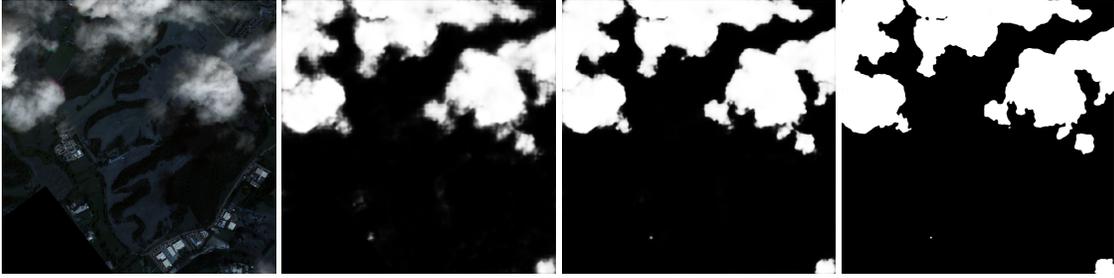


Figure 3: Examples of inferred cloud masks as determined using different loss functions for the model during training. From left to right: The inference RGB image; the model’s inferred cloud mask trained with BCE loss, trained with a loss that is the average of the BCE and Jaccard losses, trained with Jaccard loss.

### 3.3 SW-CS Dataset

The results of our model on 399 SW-CS dataset inference images is shown in table 3.3. The inference dataset contains images from various satellites with different spatial resolutions, and our model is robust enough to achieve a BOA of 94.7% and an F1 score of 92.6% even in this situation. A random selection of inference results from the model are shown in Figure 4. The model capably learns to identify opaque clouds but struggles with thin clouds, as seen in the fourth row of 4 (though again labels for thin clouds were not included for the model during training). Even with the relatively large input image sizes of  $512 \times 512$  pixels, the time to inference is  $< 2$  seconds on CPU (3.3-GHz Intel Xeon Scalable processor) and a mere 190-ms on a NVIDIA K80 GPU with 2,496 parallel processing cores and 12-GB of GPU memory.

OA	BOA	PA	UA	F1 <sub>total</sub>	IoU <sub>total</sub>	R <sub>mean</sub>	P <sub>mean</sub>	F1 <sub>mean</sub>	IoU <sub>mean</sub>
95.8	94.7	92.1	93.0	92.6	86.1	85.6	85.5	81.4	74.5

Table 2: Our Attention ResUNet’s performance trained and tested on the SW-CS data.

### 3.4 Cloud Reference Dataset

Following the results presented in the recent Cloud Mask Inter-comparison eXercise [Skakun et al., 2022] we test our model using existing L8 and S2 reference cloud datasets, which include the datasets of Hollstein [Hollstein et al., 2016], PixBox [Paperin, 2021a,b], L8Biome [Foga et al., 2017] GSFC [Skakun et al., 2021], and CESBIO [Baetens and Hagolle, 2018]. This allows us to evaluate the performance of our model using common points of comparison against which the most popularly employed cloud detection algorithms have also been evaluated [Zhu and Woodcock, 2012, Chen et al., 2021, Hu et al., 2021]. In each inter-comparison exercise we compare our Attention ResUNet against the the best performing model on each of those datasets. Since our model is trained to detect opaque clouds, it is tested on the opaque cloud annotations of these different datasets when such classes are available. No distinction is made between thick and thin clouds in the case of the GSFC-L8 and CESBIO datasets.

#### 3.4.1 Inter-class Data Pre and Post Processings

The following pre-processing steps are taken in order to prepare the image/label pairs of each external dataset. Images are cropped to the minimum area containing all annotated clouds for the Hollstein and GFSC datasets. All images are zero-padded and tiled into  $512 \times 512$ -pixel images. Labels are created such that a pixel value of 0 is transparent, 127 represents no data (and is not considered in quantitative performance evaluations), and a pixel value of 255 is denotes a cloud pixel. Labels are padded to meet the padding of the corresponding input image. No pre-processing steps for changing the contrast or spatial resolutions of the images are taken.

Finally, the output tiles (cloud masks) produced on inference by the model are combined to cover the full original image. Only in the case of the CESBIO dataset is the full image resized to fit the ground truth label dimensions because the ground truth labels for that dataset are originally provided at a smaller size than the input images.

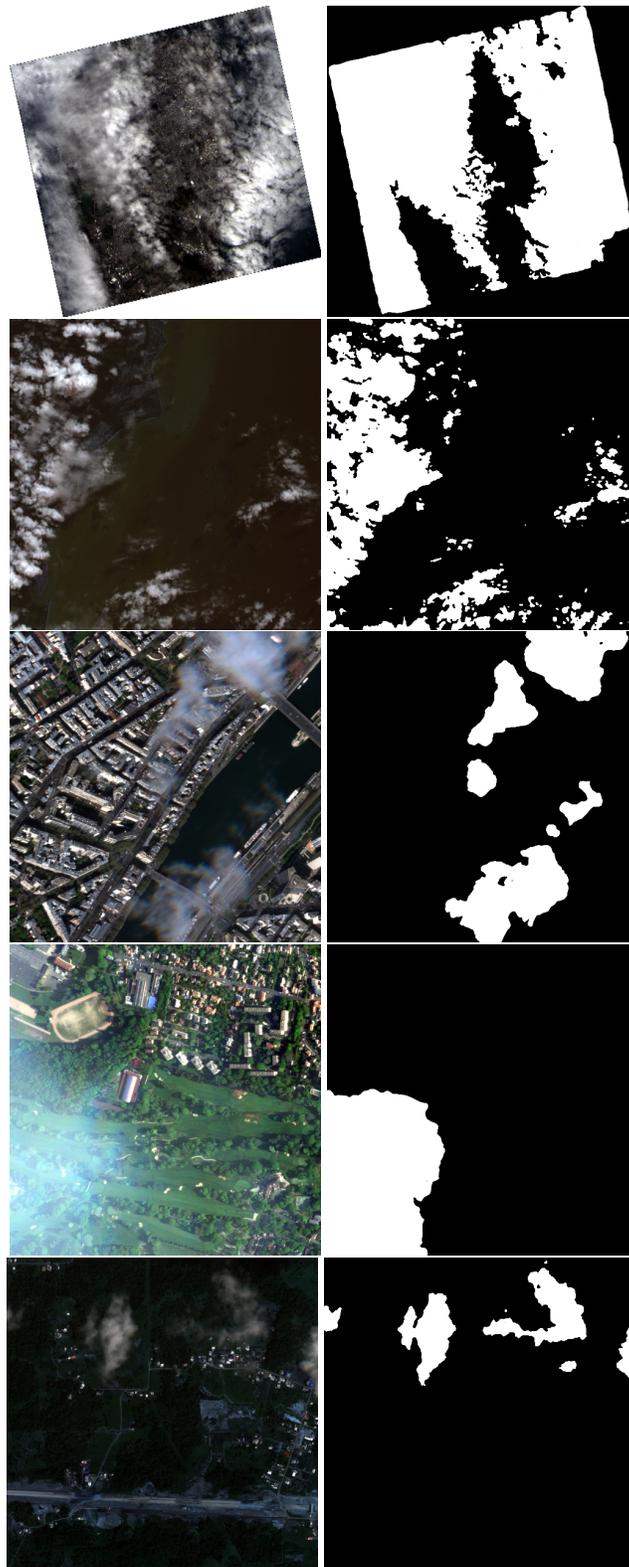


Figure 4: Examples of inference results by the model on the SW-CS dataset. The columns are: The input RGB image (left), and the model's predicted cloud mask (right).

### 3.4.2 Hollstein

The Hollstein dataset includes 59 scenes obtained by S2 (at 20-m pixel resolution). The cloud labels for this dataset are provided as polygons, meaning only some of the pixels in each scene are actually annotated explicitly. Table 3 shows the quantitative inference results of our model on the Hollstein dataset (without thin clouds). From the CMIX the CD-FCNN model was the best performing model. CD-FCNN is itself a deep learning model that also uses a UNet-based architecture and is trained on L8Biome and SPARCS datasets with a spatial resolution of 30-m to identify clouds. Our model performs comparably to the CD-FCNN on this dataset even without having been exposed to images at this resolution during training.

Algorithm	OA	BOA	PA	UA
CD-FCNN	<b>97.8</b>	<b>97.8</b>	98.3	<b>96.7</b>
Ours	96.9	97.2	<b>99.7</b>	93.7

Table 3: Hollstein - Opaque clouds only.

### 3.4.3 GSFC

GSFC datasets consists of two parts. The first part includes 6 scenes of L8 data, namely GSFC-L8. The GSFC-L8 dataset does not distinguish between thin and opaque clouds in their annotations. Table 4 shows the performance of our model compared to the top performing algorithms mentioned in [Skakun et al., 2022]. Although CD-FCNN takes a similar approach to our own and is specifically trained on this dataset, our model is able to achieves the same level of performance without having ever seen the data during training.

The second part of the dataset is referred to as GSFC-S2 and is comprised of 28 scenes of S2 data. The clouds are annotated with polygon shapes, with separate annotations for translucent and opaque clouds. Table 5 shows the performance metrics of our model and the other high-ranked algorithms. The results show that our model achieved 100% UA, which means there were no false positives in the predictions. However, the PA is lower than other algorithms indicative of the high number of false negative predictions made by the model. The reason is that there is a discrepancy between GSFC-S2 cloud annotations and what our model considers to be cloudy. In fact, the ground truth of GSFC-S2 is not precise. Many pixels are incorrectly classified as clouds when they clearly belong to the background class. In other words, the ground truth cloud masks significantly over-classify cloudy regions in these images. In spite of our model’s apparent performance, Figure 5 shows this discrepancy between the RGB images, the ground truth cloud masks, and our model’s substantially more accurate cloud masks.

Algorithm	OA	BOA	PA	UA
CD-FCNN	97.3	97.3	94.6	100.0
Ours	97.3	97.3	94.6	100.0

Table 4: GSFC-L8 - No distinction made between translucent versus opaque clouds, and for which many pixels are incorrectly classified in the ground-truth labels.

Algorithm	OA	BOA	PA	UA
LaSRC	<b>98.0</b>	<b>97.9</b>	<b>98.5</b>	97.8
CD-FCNN	92.9	93.6	87.3	99.9
Ours	85.4	86.6	73.2	<b>100.0</b>

Table 5: GSFC-S2 - Opaque clouds only. Many pixels are incorrectly classified in the ground-truth labels of this set as well.

### 3.4.4 L8Biome

L8Biome is another dataset used in [Skakun et al., 2022] which includes 96 scenes of L8 data with a spatial resolution of 30-m. In this dataset regions with thin translucent clouds belong to a separate class and are distinguished from more opaque clouds. Every image is fully annotated at the pixel level. Table 6 shows the inference results of the best algorithms mentioned in [Skakun et al., 2022] as well for our model. These models are applied on 80 out of the 96 scenes. Our model has a similar performance to the top performing algorithm on this dataset. Note that the CD-FCNN is not measured against this dataset as the L8Biome dataset is a subset of the data that model was trained on.

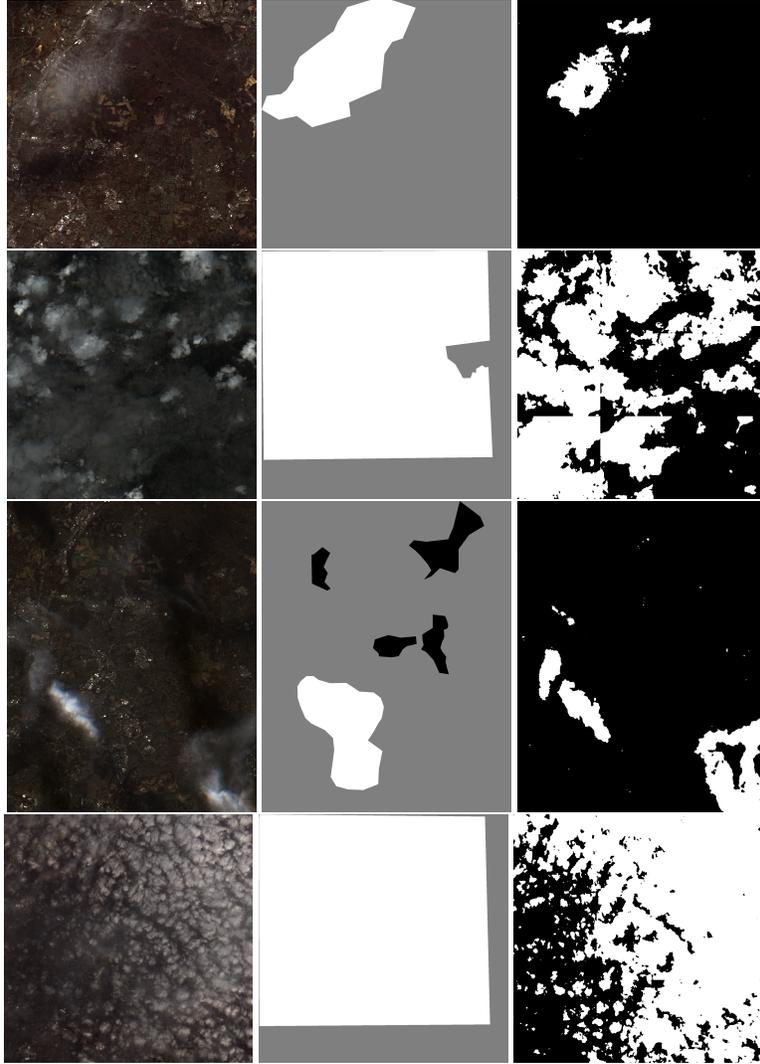


Figure 5: Examples of label discrepancies in the GSFC-S2 dataset between the true color image (left), and the ground truth cloud mask (center), compared to the inferred cloud mask predicted by our model (right). This type of misclassifications is common throughout all 28 samples in this dataset.

Algorithm	OA	BOA	PA	UA
LaSRC	<b>92.8</b>	<b>93.5</b>	97.8	<b>86.9</b>
Ours	92.1	92.9	<b>98.6</b>	85.5

Table 6: L8Biome - Opaque clouds only.

### 3.4.5 CESBIO

The CESBIO dataset is composed of S2 imagery. It consists of 30 scenes with a spatial resolution of 60-m. Each image is fully annotated at the pixel level, and there is no distinction made between thin versus opaque clouds. Table 7 shows the performance of our model and other top-performing algorithms on this dataset from the CMIX study. Our model performs worse in comparison to the CD-FCNN deep learning model, particularly on the PA and UA metrics. Our model produces a large number of false negatives because of discrepancies between the ground truth labels in our proprietary training set. In the SW-CS dataset we consider translucent clouds as belong to the "background" class in the ground truth annotations. A few examples of this discrepancy are shown in Figure 6 in which our model detects opaque clouds well, but ignores thin translucent clouds. In Table 7, a relatively low UA means that there are also a large number of false positives in our model's results. This shortcoming is explained by examining Figure 7. Not only

does our training data lack S2 samples and or images with a spatial resolution at 60-m, it also does not include enough samples of textured arid regions as shown in Figure 7. The similarity between the texture of arid scenes and of clouds thus confuse the model in the absence of sufficiently similar training data.

Algorithm	OA	BOA	PA	UA
CD-FCNN	89.5	79.5	60.3	<b>94.1</b>
FORCE	<b>91.1</b>	<b>88.9</b>	<b>84.7</b>	79.9
Ours	87.5	79.2	63.0	81.5

Table 7: CESBIO

### 3.4.6 PixBox

PixBox consists of two datasets. Both datasets are partially annotated, meaning only selected pixels in each image have been classified. In both datasets, thin clouds are included in a different class than opaque clouds. The first dataset includes 29 scenes of S2 data at 10-m spatial resolution. The inference results on PixBox-S2 without thin clouds are provided in table 8. Our model outperforms every other model in this category, achieving an OA of 95.8% and BOA of 96.5%. The second part of the dataset includes L8 imagery with a spatial resolution of 30-m in 11 scenes. Table 9 shows the inference results where our model again outperforms every other model in this category as well, as found in Skakun et al. [2022]; achieving an OA of 98.6% and BOA of 98.9%.

Algorithm	OA	BOA	PA	UA
InterSSIM	91.9	90.7	86.2	<b>91.3</b>
S2cloudless	91.6	91.6	91.6	86.4
CD-FCNN	89.5	88.1	82.7	87.9
Ours	<b>95.8</b>	<b>96.5</b>	<b>97.6</b>	84.1

Table 8: PixBox-S2 - Opaque clouds only.

Algorithm	OA	BOA	PA	UA
ATCOR	98.4	96.7	94.1	<b>95.6</b>
CD-FCNN	97.8	98.7	<b>99.9</b>	87.4
Ours	<b>98.6</b>	<b>98.9</b>	99.4	93.0

Table 9: PixBox-L8 - Opaque clouds only.

## 3.5 Error Sources

Most of the errors in the inference dataset are due to haze and small clouds not being detected by our model. The magnitude of these errors is tolerable for two key reasons. 1. The network was not trained to detect thin clouds such as haze, and 2. Small clouds do not contribute significantly to the overall cloud cover percentage.

However, there are some problematic inference samples that do not belong to these two types. Most of the errors are caused by patterns similar to clouds (e.g., patterns in arid areas, structures, and some highly textured snow scenes). Figures 8 and 7 show a few examples of these problematic inference images, from our own inference dataset and from the CESBIO dataset, respectively. The first row of images in Figure 8 show an example of farm land with a repetitive pattern of crops that our model incorrectly classifies as clouds. In the second row several buildings are classified as clouds as well. More samples (and augmentation) of these scenes can be added to the training dataset to help the network learn more about them. Our model accurately distinguishes cloud from snow in a variety of scenarios. However, in scenes dominated by snow the network struggles to distinguish cloud from snow as accurately as in scenes where snow occupies only a portion of the total scene. Figure 9 shows several examples of how our model performs on different scenes containing various degrees of land and snow cover.

All of the samples with an F1 score of less than 80% were manually examined in order to determine on which samples our Attention ResUNet had the greatest difficulty in identifying clouds. We find that where the F1 score is below 80% (in 68 out of the 399 images on which the model was tested) there are greater variations in the correspondence of the ground truth cloud mask and the actual RGB image. This is evident in standard deviation of IoU scores between these two sub-samples (below and above an F1 score of 80%): The IoU of our model when the F1 score of the prediction is above 80% is  $0.90 \pm 0.05$ ; on images where the F1 score of the prediction is below 80% the IoU falls dramatically

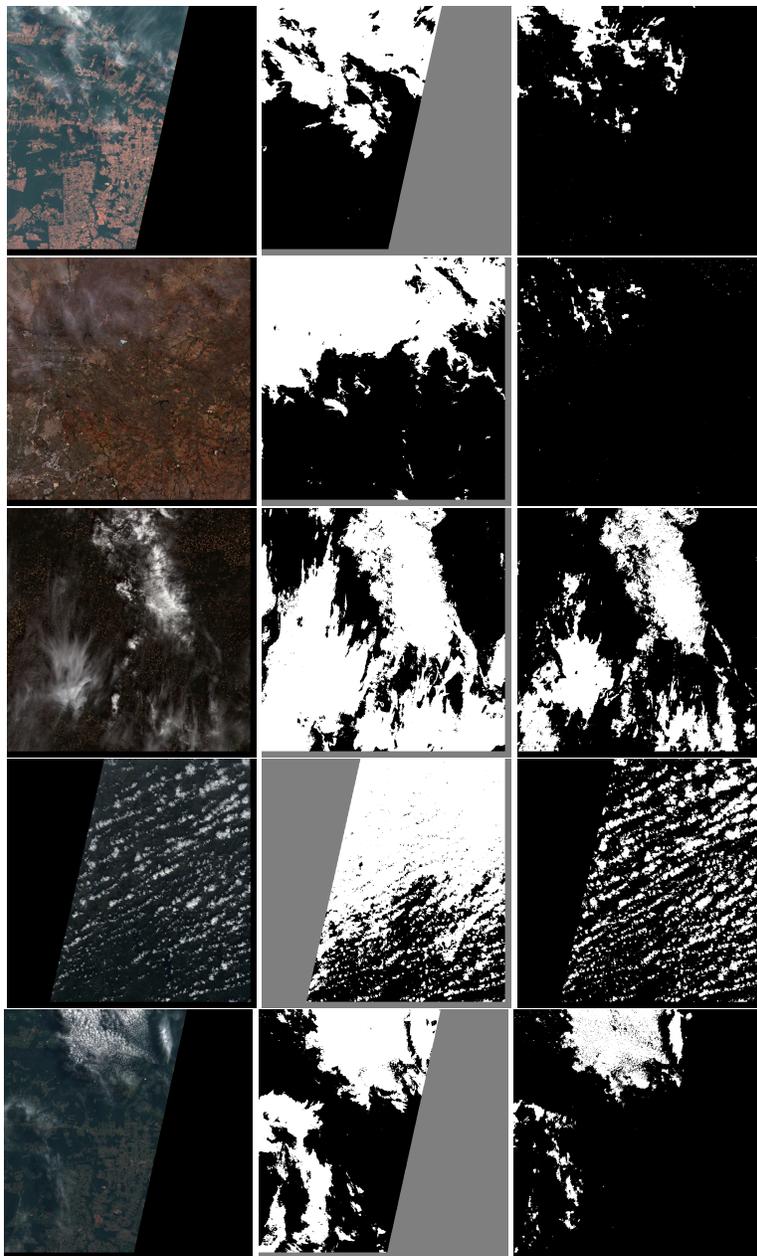


Figure 6: Examples of label discrepancies in the CESBIO dataset. Column are: True color images (left), ground truth labels (center), and our model’s predicted cloud masks (right).

to  $0.49 \pm 0.16$ , with a significant increase in the standard deviation as well. This suggests some of the labels in the ground truth of the SW-CS dataset itself may need to be revised.

#### 4 Conclusion

In this paper we have demonstrated a semantic segmentation network that is an excellent model for identifying opaque clouds in satellite imagery. Our model uses just three bands (RGB) for input and accepts relatively large image sizes ( $512 \times 512$  pixels), without the need for additional information such as from NIR bands or secondary sensor information. Our model makes use of the U-Net architecture, skip connections, and an attention mechanism. The incorporation of residual units in the encoder stage allows us to ease the training of a model of this size. Also unique to other deep learning approaches that have been employed for the semantic segmentation of clouds, is our use of

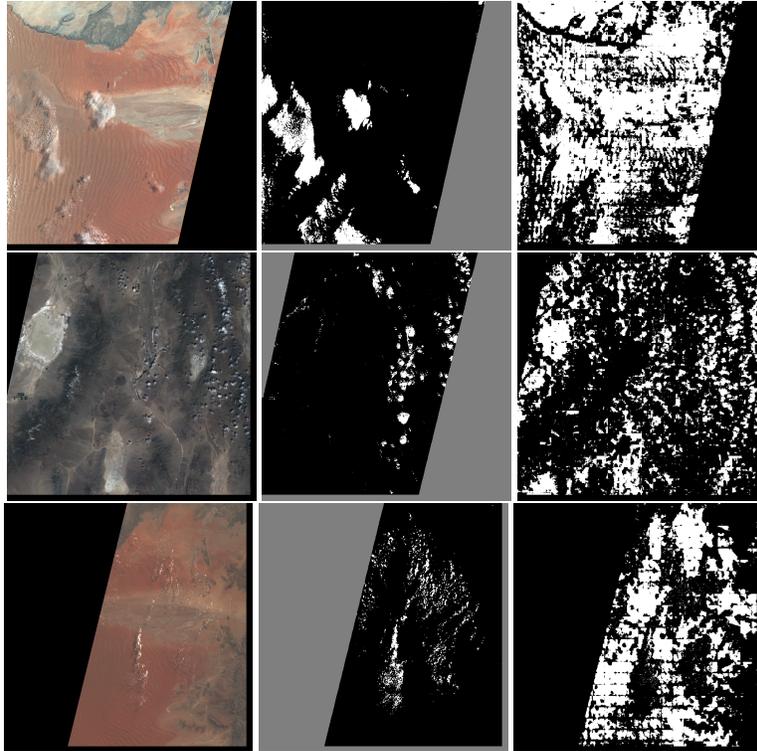


Figure 7: Examples of textured arid regions in the CESBIO dataset that our model was confused by. Columns: True color images (left), ground truth labels (center), and our model’s predicted cloud masks (right).

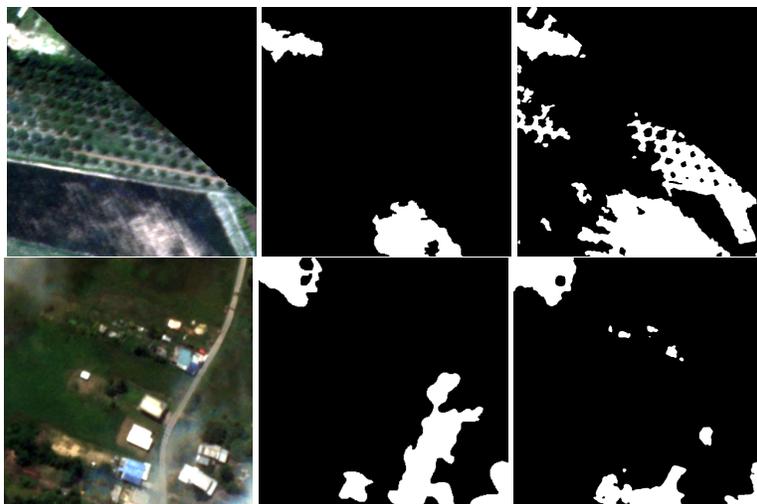


Figure 8: Examples of failure cases in our inference dataset. Columns: True color images (left), ground truth cloud masks (center), our model’s predicted cloud masks (right).

Jaccard loss during training to account for class imbalance in cases where clouds occupy considerably fewer numbers of pixels compared to the total image size. Further, our model is trained on imagery produced by multiple satellite sensors, and displays considerable robustness to changes in spatial resolution.

We evaluated our model using two experiments. The first experiment, which we refer to as an inter-class experiment, was performed on a challenging proprietary dataset. In this experiment, our model achieves an F1 score of 92.6% and a BOA of 94.7%. For the intra-class experiment we applied our model to several popular cloud segmentation datasets mentioned in Skakun et al. [2022]. Our model performed comparably to the best-in-class models tested in that paper,

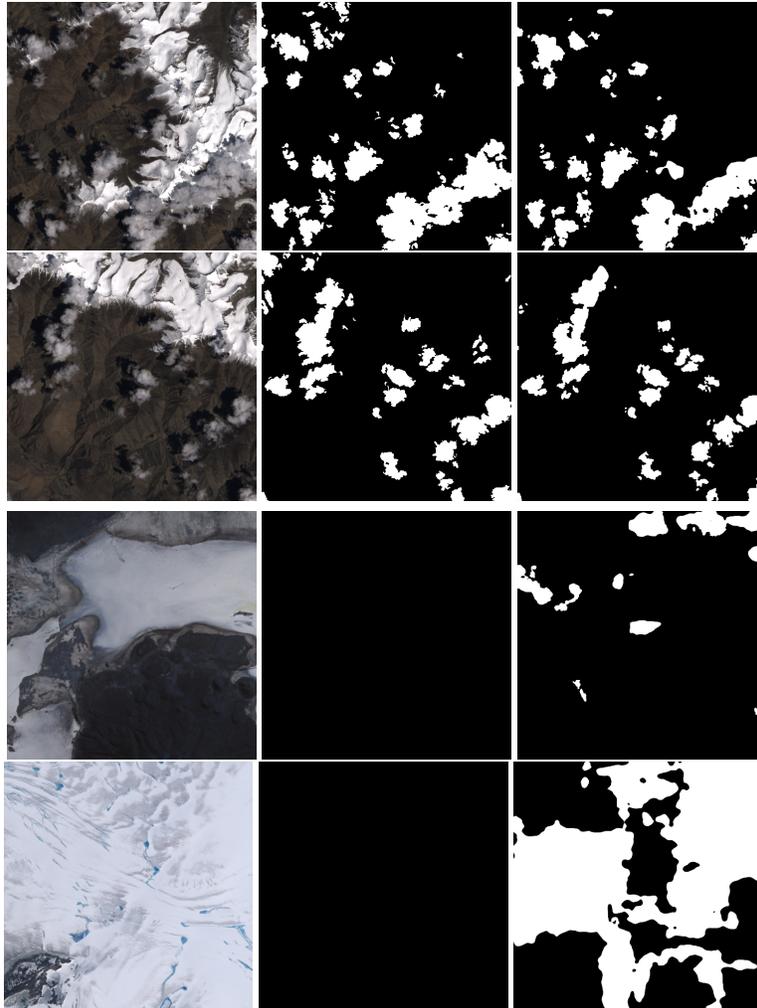


Figure 9: Examples of our model’s inferences on cloud vs snow. The first two rows show acceptable performance, and the last two rows show snowy areas misclassified as containing clouds. From left to right columns: RGB image, ground truth, our model’s prediction.

even without having been exposed to any of the L8, GSFC-L8, and L8Biome datasets, on which the models compared in that work were all trained and tested on. On both the GSFC-S2 and CESBIO datasets our model underperformed the best-in-class models, but not unexpectedly, as the ground truth labels of these datasets include a large number of false positives, which we have provided examples of herein. On the CESBIO dataset however, we found additional samples that our model found confusing in arid regions and snowy regions in which the landscape itself can appear cloud-like. On the PixBox datasets our model outperforms every other model that was tested on this dataset in Skakun et al. [2022]. Our Attention ResUNet achieves state-of-the-art performance with BOAs of 96.5% and 98.9% on PixBox-L8 and -S2, respectively. Our Attention ResUNet is therefore a novel architecture for cloud segmentation, which together with our proprietary dataset, is a significant improvement in the application of generalized cloud detection. In future studies we expect this architecture can be improved further. Knowledge distillation can reduce the number of trainable parameters of the model enabling it to be deployed on resource-limited architectures such as those on-board a satellite. We further expect General Adversarial Network (GAN) augmentation (e.g., Nyborg and Assent, 2021) could also be used to improve the model’s performance in identifying clouds over arid and snowy regions.

## References

Michael D. King, Steven Platnick, W. Paul Menzel, Steven A. Ackerman, and Paul A. Hubanks. Spatial and temporal distribution of clouds observed by modis onboard the terra and aqua satellites. *IEEE Transactions on Geoscience*

- and *Remote Sensing*, 51:3826–3852, 2013.
- David P Roy, Junchang Ju, Kristi Kline, Pasquale L Scaramuzza, Valeriy Kovalsky, Matthew Hansen, Thomas R Loveland, Eric Vermote, and Chunsun Zhang. Web-enabled landsat data (weld): Landsat etm+ composited mosaics of the conterminous united states. *Remote Sensing of Environment*, 114(1):35–49, 2010.
- Eric F Vermote, Nazmi Z El Saleous, and Christopher O Justice. Atmospheric correction of modis data in the visible to middle infrared: first results. *Remote Sensing of Environment*, 83(1-2):97–111, 2002.
- Alfredo Huete, Kamel Didan, Tomoaki Miura, E Patricia Rodriguez, Xiang Gao, and Laerte G Ferreira. Overview of the radiometric and biophysical performance of the modis vegetation indices. *Remote sensing of environment*, 83(1-2):195–213, 2002.
- Qiaofeng Zhang, J Wang, X Peng, P Gong, and P Shi. Urban built-up land change detection with road density and spectral information from multi-temporal landsat tm data. *International journal of remote sensing*, 23(15):3057–3078, 2002.
- Zhe Zhu and Curtis E Woodcock. Continuous change detection and classification of land cover using all available landsat data. *Remote sensing of Environment*, 144:152–171, 2014.
- Zhe Zhu and Curtis E Woodcock. Object-based cloud and cloud shadow detection in landsat imagery. *Remote sensing of environment*, 118:83–94, 2012.
- Zhe Zhu, Shixiong Wang, and Curtis E Woodcock. Improvement and expansion of the fmask algorithm: Cloud, cloud shadow, and snow detection for landsats 4–7, 8, and sentinel 2 images. *Remote sensing of Environment*, 159:269–277, 2015.
- Shi Qiu, Zhe Zhu, and Binbin He. Fmask 4.0: Improved cloud and cloud shadow detection in landsats 4–8 and sentinel-2 imagery. *Remote Sensing of Environment*, 231:111205, 2019.
- Zhiwei Li, Huanfeng Shen, Huifang Li, Guisong Xia, Paolo Gamba, and Liangpei Zhang. Multi-feature combined cloud and cloud shadow detection in gaofen-1 wide field of view imagery. *Remote sensing of environment*, 191:342–358, 2017.
- Dan López-Puigdollers, Gonzalo Mateo-García, and Luis Gómez-Chova. Benchmarking deep learning models for cloud detection in landsat-8 and sentinel-2 images. *Remote Sensing*, 13(5):992, 2021.
- Jiaqiang Zhang, Xiaoyan Li, Liyuan Li, Pengcheng Sun, Xiaofeng Su, Tingliang Hu, and Fansheng Chen. Lightweight u-net for cloud detection of visible and thermal infrared remote sensing images. *Optical and Quantum Electronics*, 52(9):1–14, 2020.
- Yang Chen, Qihao Weng, Luliang Tang, Qinhuo Liu, and Rongshuang Fan. An automatic cloud detection neural network for high-resolution remote sensing imagery with cloud–snow coexistence. *IEEE Geoscience and Remote Sensing Letters*, 19:1–5, 2021.
- Yongjie Zhan, Jian Wang, Jianping Shi, Guangliang Cheng, Lele Yao, and Weidong Sun. Distinguishing cloud and snow in satellite images via deep convolutional network. *IEEE geoscience and remote sensing letters*, 14(10):1785–1789, 2017.
- Min Xia, Wan’an Liu, Bicheng Shi, Liguang Weng, and Jia Liu. Cloud/snow recognition for multispectral satellite imagery based on a multidimensional deep residual network. *International journal of remote sensing*, 40(1):156–170, 2019.
- Jingyu Yang, Jianhua Guo, Huanjing Yue, Zhiheng Liu, Haofeng Hu, and Kun Li. Cdnet: Cnn-based cloud detection for remote sensing imagery. *IEEE Transactions on Geoscience and Remote Sensing*, 57(8):6195–6211, 2019.
- Kai Hu, Dongsheng Zhang, and Min Xia. Cdunet: Cloud detection unet for remote sensing imagery. *Remote Sensing*, 13(22):4533, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Ozan Oktay, Jo Schlemper, Loic Le Folgoc, Matthew Lee, Mattias Heinrich, Kazunari Misawa, Kensaku Mori, Steven McDonagh, Nils Y Hammerla, Bernhard Kainz, et al. Attention u-net: Learning where to look for the pancreas. *arXiv preprint arXiv:1804.03999*, 2018.
- Yanan Guo, Xiaoqun Cao, Bainian Liu, and Mei Gao. Cloud detection for satellite imagery using attention-based u-net convolutional neural network. *Symmetry*, 12(6):1056, 2020.
- Sorour Mohajerani and Parvaneh Saeedi. Cloud-net: An end-to-end cloud detection algorithm for landsat 8 imagery. In *IGARSS 2019-2019 IEEE International Geoscience and Remote Sensing Symposium*, pages 1029–1032. IEEE, 2019.

- Md Atiqur Rahman and Yang Wang. Optimizing intersection-over-union in deep neural networks for image segmentation. In *International symposium on visual computing*, pages 234–244. Springer, 2016.
- M Joseph Hughes and Daniel J Hayes. Automated detection of cloud and cloud shadow in single-date landsat imagery using neural networks and spatial post-processing. *Remote Sensing*, 6(6):4907–4926, 2014.
- Steve Foga, Pat L Scaramuzza, Song Guo, Zhe Zhu, Ronald D Dilley Jr, Tim Beckmann, Gail L Schmidt, John L Dwyer, M Joseph Hughes, and Brady Laue. Cloud detection algorithm comparison and validation for operational landsat data products. *Remote sensing of environment*, 194:379–390, 2017.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- Shengyuan Piao and Jiaming Liu. Accuracy improvement of unet based on dilated convolution. In *Journal of Physics: Conference Series*, volume 1345, page 052066. IOP Publishing, 2019.
- Kaili Cao and Xiaoli Zhang. An improved res-unet model for tree species classification using airborne high-resolution images. *Remote Sensing*, 12(7):1128, 2020.
- Fabian Isensee and Klaus H Maier-Hein. Or-unet: an optimized robust residual u-net for instrument segmentation in endoscopic images. *arXiv preprint arXiv:2004.12668*, 2020.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. doi:10.1109/CVPR.2009.5206848.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pages 1310–1318. PMLR, 2013.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Rich Caruana, Steve Lawrence, and C Giles. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. *Advances in neural information processing systems*, 13, 2000.
- Katarzyna Janocha and Wojciech Marian Czarnecki. On loss functions for deep neural networks in classification. *arXiv preprint arXiv:1702.05659*, 2017.
- Sergii Skakun, Jan Wevers, Carsten Brockmann, Georgia Doxani, Matej Aleksandrov, Matej Batič, David Frantz, Ferran Gascon, Luis Gómez-Chova, Olivier Hagolle, Dan López-Puigdollers, Jérôme Louis, Matic Lubej, Gonzalo Mateo-García, Julien Osman, Devis Peressutti, Bringfried Pflug, Jernej Puc, Rudolf Richter, Jean-Claude Roger, Pat Scaramuzza, Eric Vermote, Nejc Vesel, Anže Zupanc, and Lojze Žust. Cloud mask intercomparison exercise (cmix): An evaluation of cloud masking algorithms for landsat 8 and sentinel-2. *Remote Sensing of Environment*, 274:112990, 2022.
- Kay Henning Brodersen, Cheng Soon Ong, Klaas Enno Stephan, and Joachim M Buhmann. The balanced accuracy and its posterior distribution. In *2010 20th international conference on pattern recognition*, pages 3121–3124. IEEE, 2010.
- André Hollstein, Karl Segl, Luis Guanter, Maximilian Brell, and Marta Enesco. Ready-to-use methods for the detection of clouds, cirrus, snow, shadow, water and clear sky pixels in sentinel-2 msi images. *Remote Sensing*, 8(8):666, 2016.
- Wevers Jan Stelzer Kerstin Brockmann Carsten Paperin, Michael. Pixbox sentinel-2 pixel collection for cmix (version 1.0) [data set]. zenodo. <https://doi.org/10.5281/zenodo.5036991>, 2021a.
- Stelzer Kerstin Lebreton Carole Brockmann Carsten Wevers Paperin, Michael. Pixbox landsat 8 pixel collection for cmix (version 1.0) [data set]. zenodo. <https://doi.org/10.5281/zenodo.5040271>, 2021b.
- Sergii Skakun, Eric F Vermote, Andres Eduardo Santamaria Artigas, William H Rountree, and Jean-Claude Roger. An experimental sky-image-derived cloud validation dataset for sentinel-2 and landsat 8 satellites over nasa gsfc. *International Journal of Applied Earth Observation and Geoinformation*, 95:102253, 2021.
- Louis Baetens and Olivier Hagolle. Sentinel-2 reference cloud masks generated by an active learning method, October 2018. URL <https://doi.org/10.5281/zenodo.1460961>.
- Joachim Nyborg and Ira Assent. Weakly-supervised cloud detection with fixed-point gans. *CoRR*, abs/2111.11879, 2021. URL <https://arxiv.org/abs/2111.11879>.