

A geographic information-driven method and a new large scale dataset for remote sensing cloud/snow detection

Xi Wu ^{a,b,c}, Zhenwei Shi ^{a,b,c,*}, Zhengxia Zou ^d

^a Image Processing Center, School of Astronautics, Beihang University, Beijing 100191, PR China

^b Beijing Key Laboratory of Digital Media, Beihang University, Beijing 100191, PR China

^c State Key Laboratory of Virtual Reality Technology and Systems, School of Astronautics, Beihang University, Beijing 100191, PR China

^d Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI 48109, USA



ARTICLE INFO

Keywords:

Geographic information
Cloud and snow detection
Deep convolutional neural networks
Remote sensing image

ABSTRACT

Geographic information such as the altitude, latitude, and longitude are common but fundamental meta-records in remote sensing image products. In this paper, it is shown that such a group of records provides important priors for cloud and snow detection in remote sensing imagery. The intuition comes from some common geographical knowledge, where many of them are important but are often overlooked. For example, it is generally known that snow is less likely to exist in low-latitude or low-altitude areas, and clouds in different geographic may have various visual appearances. Previous cloud and snow detection methods simply ignore the use of such information, and perform detection solely based on the image data (band reflectance). Due to the neglect of such priors, most of these methods are difficult to obtain satisfactory performance in complex scenarios (e.g., cloud-snow coexistence). In this paper, a novel neural network called “Geographic Information-driven Network (GeoInfoNet)” is proposed for cloud and snow detection. In addition to the use of the image data, the model integrates the geographic information at both training and detection phases. A “geographic information encoder” is specially designed, which encodes the altitude, latitude, and longitude of imagery to a set of auxiliary maps and then feeds them to the detection network. The proposed network can be trained in an end-to-end fashion with dense robust features extracted and fused. A new dataset called “Levir_CS” for cloud and snow detection is built, which contains 4,168 Gaofen-1 satellite images and corresponding geographical records, and is over 20× larger than other datasets in this field. On “Levir_CS”, experiments show that the method achieves 90.74% intersection over union of cloud and 78.26% intersection over union of snow. It outperforms other state of the art cloud and snow detection methods with a large margin. Feature visualizations also show that the method learns some important priors which is close to the common sense. The proposed dataset and the code of GeoInfoNet are available in <https://github.com/permanentCH5/GeoInfoNet>.

1. Introduction

The fast development of remote sensing technology in the past decades has helped people better understand the earth. Optical remote sensing technology, as an important branch of the remote sensing family, is of great significance to many applications, such as target detection (Zou and Shi, 2016; Lin et al., 2017a,b; Zou and Shi, 2017), scene classification (Shi et al., 2018), etc. However, the imaging process of remote sensing images is often disturbed by clouds and snow. Previous literature shows that cloud covers on average more than half of the earth's surface every day (Zhang and Xiao, 2014; An and Shi, 2015; Xie et al., 2017; Wu and Shi, 2018). In some high latitude regions, the

ground may be also covered by snow and ice all year round. On one hand, both of the above factors will greatly affect the processing and analysis of remote sensing imagery, where the cloud can be a form of occlusion (Li et al., 2014, 2019b) and the snow might increase the reflectance sharply. On the other hand, environmental studies like climate study (Bi et al., 2019) and ecological change analysis (Campbell et al., 2005; Wang et al., 2018) require cloud/snow masks but manually labeling the images is usually time-consuming and expensive (Zhan et al., 2017). Automatic cloud and snow detection provides an efficient way of producing pixel-wise cloud/snow masks and thus forms the basis of many remote sensing applications.

Geographic information such as the altitude, longitude, and latitude

* Corresponding author at: Image Processing Center, School of Astronautics, Beihang University, Beijing 100191, PR China.

E-mail addresses: xiwu1000@buaa.edu.cn (X. Wu), shizhenwei@buaa.edu.cn (Z. Shi), zhengxi@umich.edu (Z. Zou).

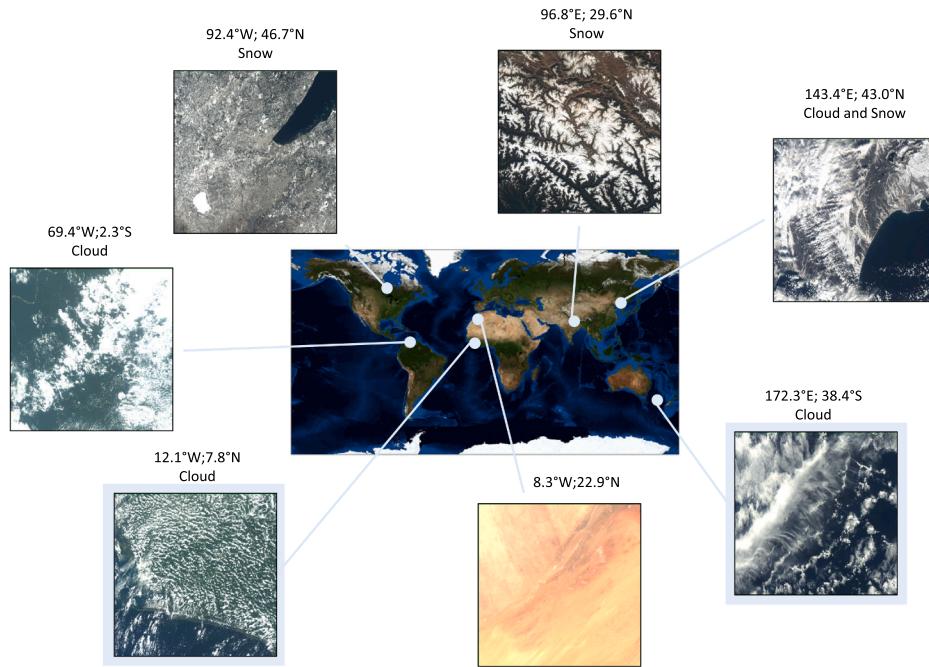


Fig. 1. Cloud and snow may represent great differences in different geographic environments. Base map credit: NASA Visible Earth.

are important meta-records in remote sensing imagery products. Such a group of records provides auxiliary and even crucial information for image processing and analysis tasks. In cloud and snow detection, it also provides important priors. For example, it is generally known that snow is less likely to exist in low-latitude or low-altitude areas, and clouds in different geographic may have various visual appearances. Fig. 1 shows some cloud and snow sample images. Each image covers around 40 thousand square kilometers, and it represents differently in different locations around the Earth. In recent years, many deep learning cloud detection and snow detection methods have been proposed. Despite the efforts made and the great improvements in this field, previous methods, even the state of the art ones, still have limitations. One of the most serious flaws of the methods is that they simply ignore the use of geographic information when performing detection. That is to say, these deep learning methods are designed solely based on the use of the image data (band reflectance), while ignoring other essential priors, such as altitude and locates. In complex scenarios such as when the cloud and snow both appear, these methods usually have difficulty generating accurate cloud and snow masks.

In this paper, a novel deep learning based method is proposed for cloud and snow detection. The method is called as Geographic Information-driven Neural Networks (GeoInfoNet). Different from the previous methods that simply focus on using image data (band reflectance) while ignoring geographic information, in the method, a “geographic information encoder” is designed, which encodes the altitude, latitude, and longitude of an image into a set of 2D maps. These maps are then integrated pixel-wisely to the detection networks and then train the whole detection model in an end-to-end fashion. It can be observed that the consistent improvement of the cloud and snow detection accuracy with the integration of the auxiliary information. The method outperforms other state of the art cloud and snow detection methods with a large margin. In addition to the new detection framework, a large dataset is also built for cloud and snow detection, which consists of 4,168 images of the Gaofen-1 satellite and is over 20 times larger than other datasets of this field. More importantly, the dataset contains the corresponding geographic information, including the longitude, latitude, and the high-resolution altitude map of each image. The contributions of this paper are summarized as follows:

1. Different from previous cloud and snow detection methods that are build based on band reflectance and simply ignore the geographic information of the imagery, a novel deep learning framework called “GeoInfoNet” is proposed which integrates the geographic information to the detection flow and learns the detection prior automatically. An encoder is designed to encode the auxiliary information such as altitude, longitude, and latitude into a set of 2D maps, which can be efficiently learned pixelwisely by the detection network in an end-to-end fashion.
2. Extensive studies on the feature visualizations are provided to show what prior knowledge the framework learns and how much the different parts contribute to the detection results.
3. A new dataset is built for cloud and snow detection, which is 20× larger than previous datasets of this task. More importantly, the geographic information along with each image is recorded in the dataset while such information is not included in previous datasets.

The following of this paper is organized as follows. In Section 2, related work to the method is introduced. In Section 3, the proposed method is introduced in detail. In Section 4, the details of the dataset Levir_CS are given. In Section 5, extensive experiments on the method are conducted and the discussions are presented in Section 6. Finally, Section 7 concludes this paper.

2. Related work

Efforts have been made for years to develop algorithms on automatic cloud and snow detection. Current methods mainly include 1) physical model based methods, 2) statistic model based methods, and 3) deep learning based methods. As for the discrimination between cloud and snow, generally, these methods are able to deal with it spectrally, spatially or temporally.

2.1. Physical model based methods

The first line of the detection methods mainly focuses on the reflectance of a specific image band or the ratio between two bands. For cloud detection, the well-known method is Automatic Cloud Cover Assessment (ACCA) (Irish, 2000; Irish et al., 2006), which is designed

based on the 2nd-6th bands of the Landsat-7 ETM+ imagery. However, this method fails on detecting warm cirrus clouds and does not produce cloud shadow masks (Zhu and Woodcock, 2012). A method named Function of masks (Fmask) (Zhu and Woodcock, 2012, 2014) is therefore proposed and can be viewed as an extension of the ACCA. The Fmask takes more bands into consideration with more physical tests, such as the whiteness test, haze optimal transformation (HOT) test, and water test. In Qiu et al. (2017), a modified version of the Fmask called Mountainous Fmask (MFmask) is proposed for cloud detection in the mountainous area. In Qiu et al. (2019), 'Fmask4.0' improves cloud detection by analyzing the spectral variability probability.

There are also some other physical model based methods proposed recently (Sun et al., 2017; Zhu and Helmer, 2018; Hagolle et al., 2010; Bian et al., 2014, 2016; Zhong et al., 2017; Li et al., 2017). In Sun et al. (2017), similar to the ACCA, the information of single band, multiband, band ratio, and band difference are extracted to detect clouds on Landsat8, NPP-VIIRS and MODIS imagery. In Zhu and Helmer (2018), different indexes such as HOT, the relative difference (RD) and the shadow index (SI), are calculated for cloud and cloud shadow detection on Sentinel-2 imagery. In Hagolle et al. (2010), cloud masks are generated by band reflectance relationships among visible and near-infrared bands of multi-temporal VENmS, FORMOSAT-2, Sentinel-2 and Landsat series imagery. Similarly, In Bian et al. (2014, 2016), multi-temporal information of spectral bands has been used in the cloud detection of HuanJing-1 satellite images. In Zhong et al. (2017), a modified ACCA algorithm is developed for discriminating the cloud from the clear background by using the 2nd-4th bands of the GaoFen-1 WVF and HuanJing-1-CCD imagery. In Li et al. (2017), similar to the Fmask (Zhu and Woodcock, 2012, 2014), relationships between all bands of the GF1-WVF imagery are considered for cloud detection.

As for the snow detection, the most commonly used method is the Normalized Difference Snow Index (NDSI) (Irish, 2000; Irish et al., 2006; Zhu and Woodcock, 2012; Hagolle et al., 2010; Selkowitz and Forster, 2016; Choi and Bindschadler, 2004). This method separates cloud and snow pixels by computing the ratio between the difference and the sum of the green band (0.52–0.60 μm) and the short-wave-infrared band (1.55–1.75 μm). The mechanics behind this method is that the snow is much more reflective in the green band than in the short-wave-infrared band (Irish, 2000; Irish et al., 2006; Zhu and Woodcock, 2012; Hagolle et al., 2010; Selkowitz and Forster, 2016; Choi and Bindschadler, 2004), therefore, a pixel with a higher NDSI indicates it is more likely to be covered by snow.

Despite the wide applications of NDSI, this method still has some limitations. In Zhu and Woodcock (2014), the authors mentioned that NDSI values of the snow covered forest areas are much lower than the pure snow pixels. To overcome this problem, a modified Norwegian Linear Reflectance-to-Snow-Cover algorithm (NLR) (Andersen, 1982) is proposed for generating better snow masks (Zhu and Woodcock, 2014). In Qiu et al. (2017), by using Digital Elevation Models (DEMs), the temperature-elevation relationship is established for better discrimination between the clouds and snow/ice pixels on the mountain area. In Wang et al. (2018), by substituting the information of green band to the near-infrared band (0.845–0.885 μm), the normalized difference forest snow index (NDFSI) is proposed for better detecting snow pixels in the forest area.

By comparing to pre-defined thresholds, the above physical model based methods can efficiently generate cloud and snow masks of the input imagery. The advantage of these methods is that they do not require any pixel-wise labels or any training process. However, these methods may also heavily rely on the reflectance of the imagery band and the pre-defined thresholds, which lacks flexibility and robustness under complex scenarios. More importantly, these methods may also fail in situations where some spectral bands, especially some short-wave-infrared bands, are not provided in the imagery (Li et al., 2017).

2.2. Statistic model based methods

To overcome the above problems, some statistic model based methods are proposed which aim to design more representative image features, along with the use of machine learning techniques. Different from those physical model based methods that only consider band features, the statistical methods use also consider spatial image features such as edge and texture, which further explore the information behind the imagery. The statistic model based methods (Li et al., 2015; Hollstein et al., 2016) usually deal with the cloud/snow detection under a classification paradigm. Some well-known classifiers such as support vector machine (SVM) (Cortes and Vapnik, 1995) and random forest (Breiman, 2001) are commonly used in this task.

The most frequent used image features for cloud and snow detection include the following types: brightness features, texture features, and local statistical features. Brightness features, i.e., the reflectance of the image bands, are the most commonly used features in statistical methods (Li et al., 2015; An and Shi, 2015). As cloud and snow pixels are usually with high reflectance, brightness features are often the first to be considered to separate the cloud/snow apart from the background. Besides, there are also methods to convert the band values to other color spaces to enrich the brightness features. For example, in Kang et al. (2018), the RGB image bands are converted to the Hue-Saturation-Intensity (HSI) color space. In Hollstein et al. (2016), the band-differences, band-ratios, and other generalized indexes, are also computed. For the local statistical features and texture features (An and Shi, 2015; Li et al., 2015; Deng et al., 2018; Srivastava and Stroeve, 2003; Kang et al., 2018), these features are usually computed by sliding windows across the whole image. By setting different window sizes, features can be extracted in multiple scales. In An and Shi (2015), Kang et al. (2018), the mean pixel value and the variance within the window are used as the local statistical feature for cloud detection. Besides, the gradient features, the Gabor filter features (Jain and Farrokhnia, 1990; Mehrotra et al., 1992; Weldon et al., 1996), and the gray level co-occurrence matrix (GLCM) (Munkres, 1957; Osada et al., 2002) are also commonly used for cloud or snow detection (An and Shi, 2015; Chen and E, 2007; Deng et al., 2018; Kang et al., 2018; Li et al., 2015; Musial et al., 2014).

2.3. Deep learning based methods

In recent years, deep neural networks (Simonyan and Zisserman, 2015; He et al., 2016; Huang et al., 2017) have made great breakthroughs in many computer vision tasks such as image classification, object detection, etc. The deep neural networks also have greatly promoted the research of cloud and snow detection in the remote sensing field.

Some early attempts of this group of methods consider the cloud and snow detection task as a patch-by-patch image classification process (into three categories: "cloud", "snow" and "background") (Shi et al., 2016; Xie et al., 2017; Le Goff et al., 2017; Zi et al., 2018; Mateo-García et al., 2017). In Mateo-García et al. (2017), cloud detection is conducted on 33×33 image patches with a 2-layer convolutional neural network. In Le Goff et al. (2017), Zi et al. (2018), Shi et al. (2016), a pixel cluster method called simple linear iterative cluster (SLIC) (Achanta et al., 2012) is first used to segment the image into a set of super-pixels, and then neural networks are used to classify each of these super-pixels. In Xie et al. (2017), a modified SLIC method is proposed and a modified AlexNet (Krizhevsky et al., 2012) is applied to further predict whether these super-pixels are covered by thick or thin cloud. The above methods take advantage of the power of deep learning neural networks in image classification and obtain higher accuracy than the statistic model based methods. However, the patch-based detection methods also have some limitations. The first limitation is that it fails when the image patch contains pixels from multiple classes and the second one is that the model only perceives locally and ignores the information from the

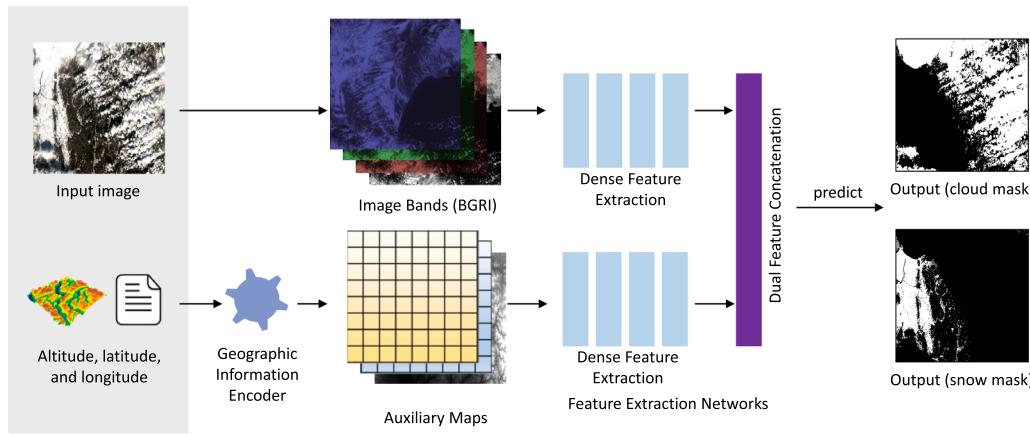


Fig. 2. Overview of the proposed method. A new network called Geo-InfoNet is proposed that makes use of both the image data and the auxiliary geographic information for cloud and snow detection. A Geographic Information Encoder is designed to encode this piece of information into a set of auxiliary maps. Features of both network-branches are extracted by “Feature Extraction Networks”, which includes “Dense Feature Extraction” module and “Dual Feature Concatenation” module. The former module can extract representative features of each branch, while the latter module is designed to produce the refined feature representation, which is further used for generating cloud and snow masks.

surrounding patches.

To overcome these limitations, the fully convolutional networks (FCNs) (Long et al., 2015) are recently introduced for cloud and snow detection (Zhan et al., 2017; Wieland et al., 2019; Jeppesen et al., 2019; Chai et al., 2019; Yan et al., 2018; Wu and Shi, 2018; Shao et al., 2019; Yang et al., 2019; Li et al., 2019a, 2020). This group of methods frames the cloud and snow detection as a pixel-wise semantic segmentation process. In Zhan et al. (2017), a VGG16-based (Simonyan and Zisserman, 2015) fully convolutional network is applied for cloud and snow detection. In Wieland et al. (2019), Jeppesen et al. (2019), Chai et al. (2019), UNet (Ronneberger et al., 2015) architecture is used to study on cloud and cloud shadow segmentation. In Yan et al. (2018), a modified residual network with pyramid pooling modules is used for cloud and cloud shadow detection. In Wu and Shi (2018), Li et al. (2019a), feature fusion of different layers are introduced to improve the details of the cloud detection result. In Shao et al. (2019), multiple bands of the Landsat-8 imagery are used as the network input, where other methods usually take the RGB bands as their inputs. In Yang et al. (2019), some network units are specifically designed for cloud detection, including the context exploitation, score map resolution preservation, and boundary refinement. In Li et al. (2020), cloud detection is processed in a deep matting framework where the images of cloud reflectance, attention mask and opacity can all be obtained.

Although these deep learning methods greatly improve the detection accuracy over traditional detection methods, the use of geographic information is not explored yet. This is one of the reasons the method is designed in this paper.

2.4. Discrimination between cloud and snow

In the process of cloud and snow detection, there are usually difficult cases where the remote sensing image contains both cloud and snow. Therefore, the discrimination between cloud and snow is very significant. Generally, the above-summarized cloud and snow detection methods can also be divided into three classes in another dimension: spectral methods, spatial methods and temporal methods.

In detail, many physical methods (Irish, 2000; Irish et al., 2006; Zhu and Woodcock, 2012, 2014; Qiu et al., 2017; Qiu et al., 2019; Li et al., 2017; Selkowitz and Forster, 2016; Choi and Bindschadler, 2004; Wang et al., 2018) tend to separate cloud and snow in the spectral domain. These spectral methods analyze the relationships between the spectral bands of the remote sensing images to classify cloud and snow. Based on the spectral band relationships, many spectral filters can be established to choose cloud and snow pixels. Although the filters of obtaining cloud masks are generally different among these methods, the approaches to generate snow masks are almost the same. Specifically, the core of the

spectral snow detection methods is NDSI, which is an index utilizing the green band and the short-wave-infrared band. With the previously defined thresholds, these spectral methods are able to produce cloud and snow masks of the input remote sensing images.

The second part is spatial methods. These methods are generally the above mentioned statistical model based methods (An and Shi, 2015; Chen and E, 2007; Deng et al., 2018; Hollstein et al., 2016; Kang et al., 2018; Li et al., 2015; Musial et al., 2014; Srivastava and Stroeve, 2003) and deep learning based methods (Shi et al., 2016; Xie et al., 2017; Le Goff et al., 2017; Zi et al., 2018; Mateo-García et al., 2017; Zhan et al., 2017; Wieland et al., 2019; Jeppesen et al., 2019; Chai et al., 2019; Yan et al., 2018; Wu and Shi, 2018; Shao et al., 2019; Yang et al., 2019; Li et al., 2019a, 2020). The spatial methods try to use the spatial information of the image by extracting manual designed image features or deep learning network features. After image features are obtained, there are usually pre-trained classifiers or classifying network layers to discriminate cloud and snow according to these spatial features. To make the classifiers robust, enough training data are needed in these spatial methods. Therefore, spatial methods can also be viewed as data-driven methods.

A few methods try to discriminate cloud and snow by using multi-temporal information of the input images (Hagolle et al., 2010; Bian et al., 2014, 2016). Among this type of methods, multi-temporal tests are introduced by using the blue band to detect the cloud pixels. The multi-temporal tests can utilize the time-series images and are very efficient to obtain cloud pixels. For snow detection, NDSI (Hagolle et al., 2010) or whiteness (Bian et al., 2014, 2016) tests are used to extract snow pixels, which is similar to spectral methods. As these temporal methods need image series, they are not proper in the situations where the input is limited to only one remote sensing image scene.

It should be noted that the discrimination between cloud and snow is a challenging task mainly due to two aspects of reasons. One is that cloud and snow represent a high visual similarity in remote sensing images, especially in visual bands. For example, both of them have high reflectance in most bands. Also, they all have irregular shapes and appear in different scales in remote sensing images. Hence, many spectral tests can not distinguish cloud and snow very well and are heavily relied on pre-defined thresholds. These methods such as (Li et al., 2017; Qiu et al., 2017; Qiu et al., 2019) may produce wrong-detection results that misclassify snow as cloud. The other reason is the imbalanced data distribution. Generally, it is more likely to see clouds in remote sensing images. On our planet, more than half of the surface is covered by clouds every day (Zhang and Xiao, 2014; An and Shi, 2015; Xie et al., 2017; Wu and Shi, 2018), while snow is seldom witnessed in some regions, such as low-latitude or low-altitude areas. Therefore, the spatial methods will more pay more attention to clouds

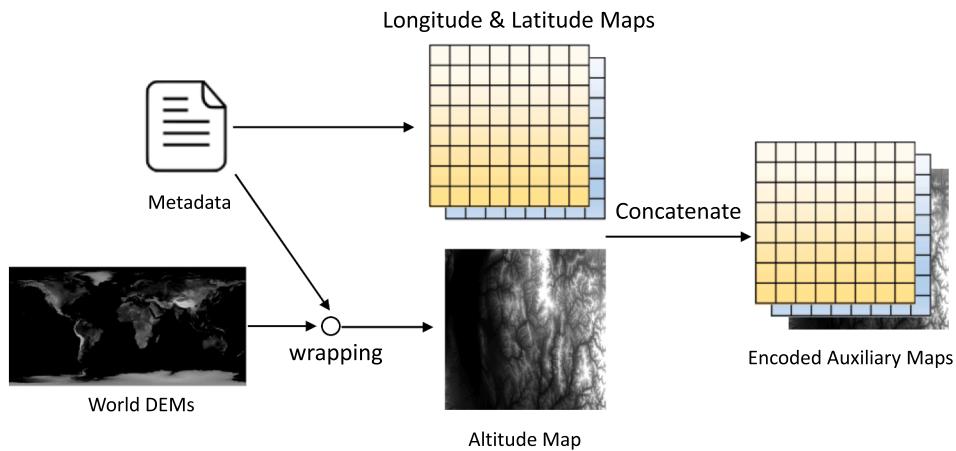


Fig. 3. The processing pipeline of the Geographic Information Encoder. For an input image, the Longitude map A_{Long} and the Latitude map A_{Lat} are generated from the metadata based on the Eq. (1). The altitude map A_{Alt} of the input image is also generated by pixelwisely wrapping the world DEMs to the image projection coordinates. Finally the three maps are concatenated to produce the final encoded auxiliary maps.

rather than snow mainly because there are much more image data with clouds than snow. As a result, the area of snow may be more likely to be detected as clouds by using the spatial methods as the visual results shown in Yang et al. (2019), Li et al. (2020). Above all, the spectral and the spatial information may not be enough in the discrimination of cloud and snow, and using geographic information can be one potential solution.

Despite there are many challenges in the discrimination between cloud and snow, there are no public datasets for snow detection, which may limit the pace of the related research. The proposed Levir_CS dataset in this paper may contribute to this research field.

3. Methodology

In this section, a detailed description of the detection method is given and how the geographic information is encoded and integrated to the network.

3.1. An overview of the GeoInfoNet

Fig. 2 shows an overview of the method. The proposed GeoInfoNet is an end-to-end network that utilizes both the input image and a set of auxiliary maps. The auxiliary maps are produced by the Geographic Information Encoder, which will be introduced in Section 3.2. In Geo-InfoNet, the network structure in DenseNet (Huang et al., 2017) is followed as the backbone network and extract multi-scale dense features from the input image and auxiliary maps separately. Then these features extracted from the two branches will be merged and used to produce the final cloud and snow masks. The two basic modules in the method, i.e., the “Dense Feature Extraction” and the “Dual Feature Concatenation”, which form “Feature Extraction Networks”, will be described in detail in Section 3.3.1 and Section 3.3.2.

3.2. Geographic information encoder

A geographic information encoder is designed to encode three types of meta-records along with the imagery, i.e., the longitude, latitude, and altitude, into a set of auxiliary maps. This module can be viewed as a pre-processing module of the proposed GeoInfoNet. These maps are generated to be the ones with the same spatial size of the input image but may have a different number of channels. Fig. 3 shows the processing pipeline of the geographic information encoder.

Given a remote sensing image with the size of $h \times w$, firstly the longitude and latitude are recorded on its the top-left corner and bottom-right corner. Then the longitude map A_{Long} and the latitude map A_{Lat} are

generated through an Affine transformation model (Warmerdam, 2008; Zhao et al., 2010). For a certain pixel in row y ($0 \leq y < h$) and column x ($0 \leq x < w$), the corresponding longitude $A_{\text{Long}}(y, x)$ and latitude $A_{\text{Lat}}(y, x)$ can be calculated as followings:

$$\begin{aligned} A_{\text{Long}}(y, x) &= A_{\text{Long}}(0, 0) + y \times r_{1,1} + x \times r_{1,2} \\ A_{\text{Lat}}(y, x) &= A_{\text{Lat}}(0, 0) + y \times r_{2,1} + x \times r_{2,2} \end{aligned} \quad (1)$$

where $A_{\text{Long}}(0, 0)$ and $A_{\text{Lat}}(0, 0)$ are the longitude and latitude value of the top-left image corner. $r_{1,1}, r_{1,2}, r_{2,1}$, and $r_{2,2}$ are the longitude/latitude resolution units on x and y directions, which can be obtained from the metafile of the imagery product or can be estimated from the coordinates of the four image corners and the center point.

In addition to the longitude and latitude, the altitude of the image is also encoded to another auxiliary map A_{Alt} . Given an image along with its corresponding longitude/latitude information, the altitude map A_{Alt} can be generated by pixel-wisely wrapping the global Digital Elevation Models (DEMs) to the projection coordinates of this image. For most optical remote sensing imagery products, the image altitude information is not included in the metafile. In the paper, the used DEMs are created based on the data collected by the 2000 Shuttle Radar Topography Mission (SRTM) and the resolution is 3 arc seconds (spatial resolution: 90 m). The data can be download from the following URL: <http://viewfinderpanoramas.org/dem3.html>.

The final encoded auxiliary maps A for each input image can be represented as a concatenation of the above three maps in the channel dimension as follows,

$$A = \text{concat}(A_{\text{Alt}}, A_{\text{Long}}, A_{\text{Lat}}), \quad (2)$$

where the dimension of the concated map A is $(h, w, 3)$.

3.3. Feature extraction networks

3.3.1. Dense feature extraction

In deep learning based cloud and snow detection methods, learning robust feature representations is crucial for the detection task. Since improving backbone of the networks is not the focus of the paper, a well-known backbone called “DenseNet” (Huang et al., 2017) is simply used, which achieves state of the art results in a variety of tasks, as the backbone network to extract high quality features from the input data arrays.

The DenseNet consists multiple dense blocks. In each block, the feature from all preceding convolutional layers are concated together. Formally, the feature maps M_{l+1} of the $(l+1)^{\text{th}}$ layer can be calculated as follows,

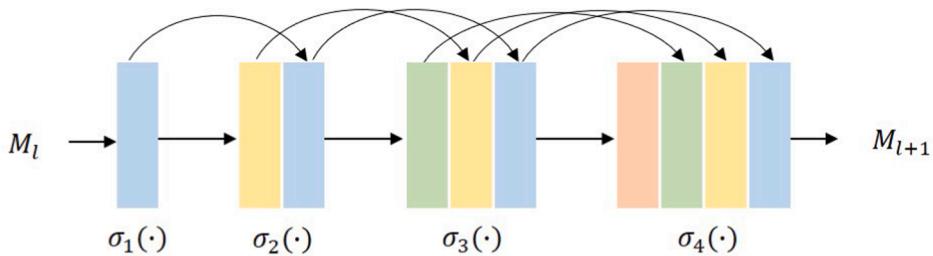


Fig. 4. An illustration of a 4-layer dense block. Each convolution layer takes all preceding feature-maps as input.

Table 1
The configuration of the dense feature extraction module.

Layers	Layer Settings
Conv_0	7×7 conv, stride=1, #out_channels=64
Pool_0	3×3 max_pool, stride=2
Dense_block_1	6 bottlenecks, #out_channels=256
Transition_1	1×1 conv, 2×2 avg_pooling, stride=2, #out_channels=128
Dense_block_2	12 bottlenecks, #out_channels=512
Transition_2	1×1 conv, 2×2 avg_pool, stride=2, #out_channels=256
Dense block_3	32 bottlenecks, #out_channels=1280
Transition_3	1×1 conv, 2×2 avg_pool, stride=2, #out_channels=640
Dense_block_4	32 bottlenecks, #out_channels=1664

$$M_{l+1} = \sigma(\text{concat}(M_l, M_{l-1}, \dots, M_1)), \quad (3)$$

where $\sigma(\cdot)$ represents a non-linear transformation on the features. Fig. 4 shows the process of the calculation of M_{l+1} .

From Eq. 3, feature maps M_1, M_2, \dots, M_l are preserved in the process of calculating M_{l+1} . Considering the concatenation of feature maps is space consuming, the filter number of each convolutional layers u is set to a small number, say, $u = 32$, compared to that in a standard convolutional network, e.g. VGG (Simonyan and Zisserman, 2015) and ResNets (He et al., 2016). In this case, the number of input feature maps in the $(l+1)^{\text{th}}$ layer will be $u_1 + 32 \times l$, where u_1 is the number of feature maps in the first layer and 32 is the filters in each layer, which also can be considered as the increasing rate. A small increasing rate not only regulates the number of features, which makes the feature extracting networks go relatively deep, but also equalizes the number of features added in each layer since the newly added information should be viewed as the same importance.

The non-linear transformation $\sigma(\cdot)$ in the networks consists of two types of operations, the normalization operation (batch normalization (Ioffe and Szegedy, 2015)), and the non-linear activation operation (rectified linear unit function (Glorot et al., 2011)). It should be noted that 1×1 convolution can be placed before 3×3 ones, which seems like

a “bottleneck”, and the settings are able to improve computational efficiency in He et al. (2016), Huang et al. (2017). Therefore, following the idea of the “Bottleneck design”, $\sigma(\cdot)$ is designed in the form of BN-ReLU-Conv(1×1)-BN-ReLU-Conv(3×3), where each Conv (1×1) outputs $4u$ feature maps.

In addition to the above dense connection module, some down-sampling modules are also designed in the networks to reduce the size of feature maps spatially and increase the computational efficiency (Wu and Shi, 2018). These modules are designed as transition blocks by following the configuration of BN-ReLU-Conv(1×1)-Pool(average, 2×2), and are placed between dense blocks. The 1×1 convolution here outputs a half number of the input feature maps.

In Huang et al. (2017), several different types of DenseNet configurations have been proposed, including DenseNet121, DenseNet169, DenseNet201 and DenseNet264. The number “X” in “DenseNetX” represents the number of convolution layers used in the classification network. In Dense Feature Extraction module, the configuration of DenseNet169 is adopted because of the balance of computation efficiency and the cost of GPU memory. The module takes in an input array which is first processed through an initial convolution layer (“Conv_0”) and an initial pooling layer (“Pool_0”), then through four dense blocks and three transition blocks accordingly. Different from the settings in Huang et al. (2017), the stride of “Conv_0” is set to 1 and remove the last classification layer for the tasks of cloud and snow detection. The configuration details of the Dense Feature Extraction module is listed in Table 1.

3.3.2. Dual feature concatenation

In the above feature extraction process, as the layers go deeper, the number of output feature maps becomes larger. The spatial resolution of the final output is down-sampled to $16 \times$ compared to that of the input, as shown in Table 1. To produce high-resolution cloud and snow masks, it is essential to increase the feature resolution by taking features. This can be done by merging the features from different blocks and generate fine-grained feature representations. The Dense Feature Concatenation module thus is designed according to this purpose. In this module, the

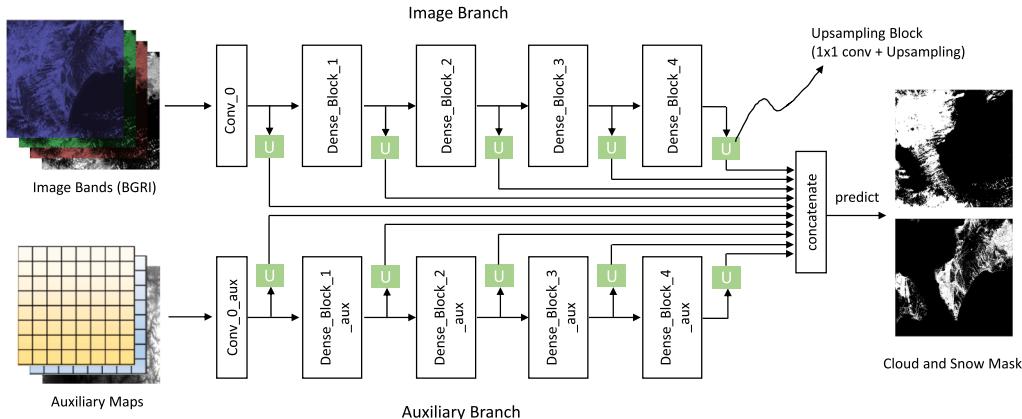


Fig. 5. Details of the Dense Feature Extraction Module. The method takes in two types of the inputs simultaneously, i.e., the original RGBI bands of the input remote sensing image and the auxiliary maps that are encoded by the Geographic Information Encoder. In both of the two branches, the dense features are extracted by the “Conv_0” layer and the following four dense blocks. The features from different blocks are upsampled to the same size and adjust the number of their channels by 1×1 convolutions, and then concatenate all these features along their channel dimension. Finally, a prediction layer is used to produce the pixel-wise score maps of the cloud and snow.

Table 2

A comparison between our dataset and the other public cloud detection datasets.

Dataset	Source	#Scenes	Snow	Geo Info
L7_Irish (Scaramuzza et al., 2011)	Landsat-7 ETM+	166	✗	✗
L8_Biome (Foga et al., 2017)	Landsat-8 OLI/TIRS	92	✗	✗
GF1_WHU (Li et al., 2017)	Gaofen-1 WVF	108	✗	✗
Levir_CS (ours)	Gaofen-1 WVF	4,168	✓	✓

initial features from both of the Blue-Green-Red-Infrared (BGRI) input image and the encoded auxiliary maps are used.

As shown in Fig. 5, for either of the two branches of the networks (i.e., input image branch and the auxiliary maps branch), the spatial features from each feature block are firstly upsampled to the size of the input image by using bilinear interpolations. Then, the upsampled features are concatenated all together along their channel dimension. Before the concatenation, we also use 1x1 convolution to adjust the channel dimension of the features from each block so that they will have the same number of channels. The intuition behind this operation is that it is assumed that for all the blocks in the networks, the features should be viewed with the same significance in the cloud and snow detection tasks. The final concatenated features M from all blocks the two branches can be represented as follows,

$$M = \text{concat}(M_{\text{img},0}, \dots, M_{\text{img},4}, M_{\text{aux},0}, \dots, M_{\text{aux},4}), \quad (4)$$

where the subscripts “img” and “aux” refer to the features from the BGRI image branch and the auxiliary information branch, respectively. The subscripts “0–4” refer to the upsampled features from the “Conv_0”, “Dense_block_1”, “Dense_block_2”, “Dense_block_3”, and “Dense_block_4” of the Dense Feature Extraction module.

3.3.3. Loss settings

In the proposed GeoInfoNet, a prediction layer (a convolutional layer with 1×1 filters) is used to produce the pixel-wise score maps of different classes: background S_1 , cloud S_2 and snow S_3 . The output score maps are normalized by using a softmax function and convert the pixel scores $(-\infty, \infty)$ to probabilities $[0,1]$. The probability map P_t of the each class $t = \{1, 2, 3\}$ can be expressed as follows,

$$P_t = \frac{\exp(S_t)}{\sum_{m=1}^3 \exp(S_m)}. \quad (5)$$

As the detection of cloud and snow is essentially a pixel-wise classification process, the networks are trained by using a standard pixel-wise classification loss (a.k.a., the cross-entropy loss). Suppose $y_m \in \{0, 1\}$ represents the ground truth label of the class m . The loss function of each pixel is expressed as follows:

$$L = - \sum_{m=1}^3 y_m \log(P_m). \quad (6)$$

Finally the average loss across all pixel from all images in the training set is computed as the final loss function.

4. Levir_CS: a new large scale dataset for cloud and snow detection

A large scale dataset called “Levir_CS” is built, where ‘C’ is for cloud and ‘S’ is for snow, respectively. As the name of the authors’ laboratory is “LEarning, VIision and Remote sensing laboratory”, similar to (Zou and Shi, 2017), the name of this dataset is started with “Levir”. Although there are already some public datasets on this topic released in the past, they are relatively small and do not contain geographic information. Besides, there are no previous public datasets for snow detection. Table 2 shows a comparison between our dataset and the other public cloud detection datasets (Scaramuzza et al., 2011; Foga et al., 2017; Li et al., 2017). Compared to other datasets listed in Table 2, the number of scenes in LEVIR_CS is over 20× larger than the other datasets, therefore, the proposed dataset is called “Large Scale”. The proposed Levir_CS dataset is available at <https://github.com/permanentCH5/GeoInfoNet>.

Gaofen-1 satellite (GF-1) is running at sun synchronous orbit, where the angle is 98.0506° and the average orbit height is 645 km. The revisiting time is 4 days. The descending node is 10:30 am. The radiometric resolution GF-1 Wide Field of View sensor (GF-1 WVF) is 10 bit. A GF-1 WVF scene, each 211 km wide by 192 km long, has an Instantaneous Field Of View (IFOV) of $16 \text{ m} \times 16 \text{ m}$ in all four bands. The spectral range is 450 nm to 890 nm. In detail, the spectral range of these bands are 450–520 nm (Blue Band or Band 1), 520–590 nm (Green Band or Band 2), 630–690 nm (Red Band or Band 3), 770–890 nm (Near Infra-

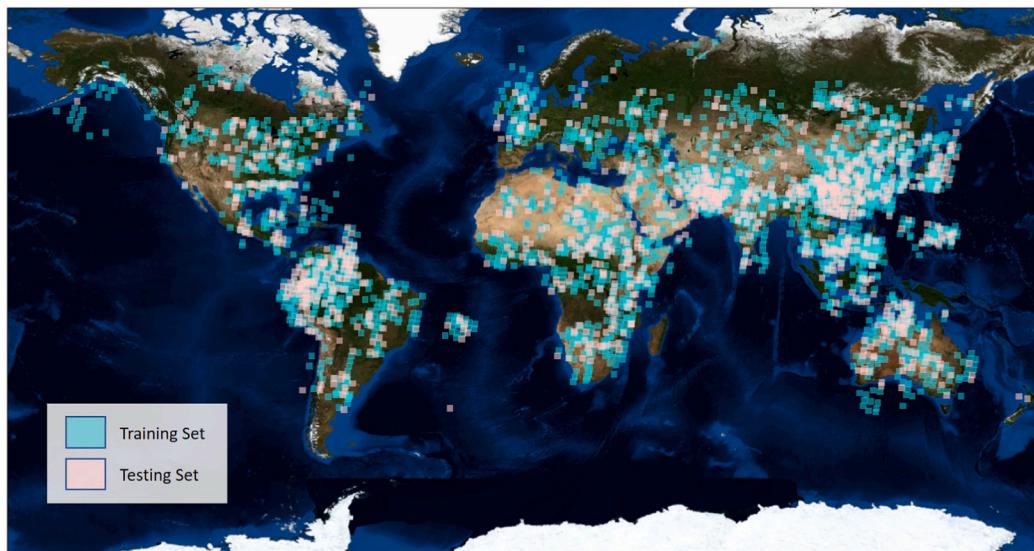


Fig. 6. The global distribution of the images in the Levir_CS dataset. Levir_CS consists of 4,168 Gaofen-1 Wide Field of View (GF-1 WVF) images. All images are obtained from the China Centre for Resources Satellite Data and Application (CRESDA) <http://www.cresda.com/>. Base map credit: NASA Visible Earth.

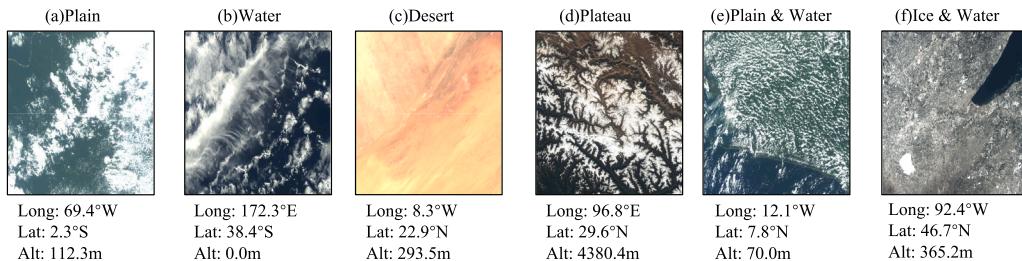


Fig. 7. Some sample scenes with different types of ground features in the proposed Levir_CS dataset. On the top of each scene, the type of the ground feature is given. On the bottom of each scene, the longitude and latitude of the central point and the mean altitude of the image is presented.

Red Band or Band 4), respectively.

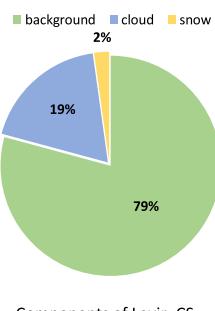
Our proposed Levir_CS consists of 4,168 GF-1 WVF scenes in total. These scenes are randomly divided into two sets, a training set with 3,068 scenes and a testing set with 1,100 scenes. The scenes in the dataset have a global distribution, as shown in Fig. 6. They cover different types of ground features, such as plain, plateau, water, desert, ice, etc. There are also combinations of the above mentioned ground feature types. Fig. 7 presents some sample scenes. Besides, as these scenes are in a global distribution, therefore, these scenes may contain different types of climate conditions, such as desert climate (see Fig. 7(c)) or sea climate (see Fig. 7(b,e)), which may help the related researches similar to (Bi et al., 2019). All the scenes were acquired from May 2013 to February 2019 and were downloaded from <http://www.cresda.com/>.

In the proposed LEVIR_CS dataset, for each scene, the level-1A product data with the process of radiation calibration is used and the current data is not produced with systematic geometric correction. This is because in many practical cases, cloud and snow detection is required to be performed in this product level to save the time of geometric correction or for fast browsing. Dataset users are able to obtain ac-

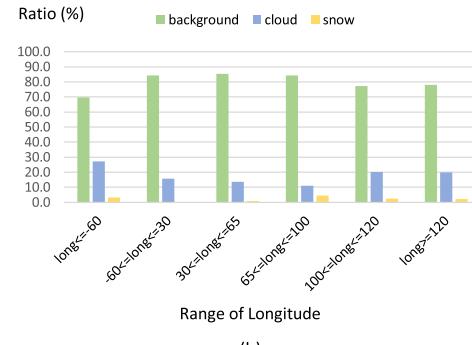
cording to the provided file of rational polynomial coefficients (RPC) to conduct systematic geometric correction if it is needed. Furthermore, to reduce the processing time of each scene and to accelerate the learning process of the global information, similar to (Zou et al., 2019), the images in LEVIR_CS dataset are 10x downsampled. For each scene in LEVIR_CS dataset, the image size is 1320×1200 and the spatial resolution is 160 m. All the four bands are used. Therefore, the resolution of DEM (90 m) is high enough in the altitude map generation. Therefore, SRTM data are chosen as the source of DEM.

In the proposed LEVIR_CS dataset, for each scene, the georeferenced multi-spectral image, the digital elevation model image and the corresponding ground truth image are all provided. The cartographic projection system used in the dataset is the World Geodetic System (WGS) and the latest version (WGS 84) is used. Through this cartographic projection system, for each scene, all the images can be registered through the geographic information. Therefore, climatic conditions do not relate to the generation of georeferenced images. For the generation of digital elevation model image, the average producing time is 45.62 s per scene.

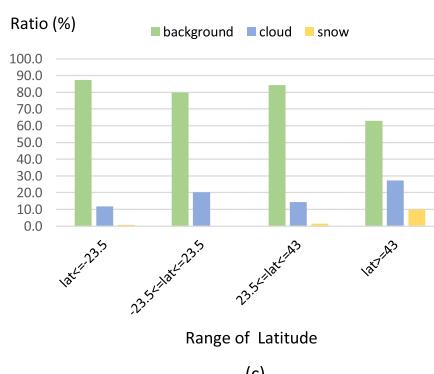
For all the images in the dataset, their pixel-wise label masks are



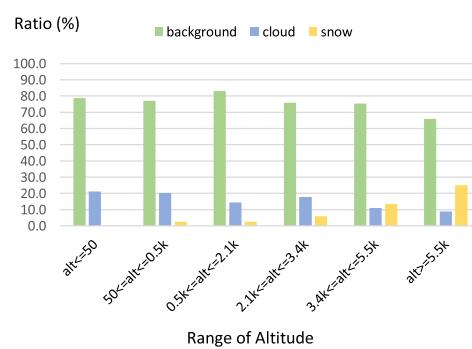
(a)



(b)



(c)



(d)

Fig. 8. Statistics of the Levir_CS dataset from different views. (a) The pixel population of the three categories: the background occupies the most (79.2%), while the snow occupies the least (2.2%). The cloud pixels occupy 18.6% of the whole pixel population. (b)-(d) Label components from longitude, latitude and altitude views, respectively. It can be seen that the distributions of the three categories are very different in different areas.

manually labeled into three categories: “background” (labeled as 0), “cloud” (labeled as 127) and “snow” (labeled as 255). Similar to (Li et al., 2017), the labeling process is finished in Adobe Photoshop. Blue, green and red bands of the original images are combined to compose a RGB image for manually labeling. To increase the labeling efficiency, similar to (Lu et al., 2019), a pre-segmentation is firstly performed by manually setting thresholds as the traditional physical methods such as (Zhong et al., 2017) indicates. Then, a rough pixel classification on these pre-segmentation region is conducted. The boundary of cloud or snow area is usually fuzzy. Like the previous research (Li et al., 2017, 2019a), these regions of the image are carefully labeled by using the brush tool (less than 10 pixels) or the lasso tool. For the thin cloud area, if the ground cover is invisible, then it is labeled as clouds. For the shadow area, as it is very dark and the region is invisible, it is labeled as the background. When labeling difficult area, the magnifying glass tool is used (more than 200% local area enlarged), which helps the labeling man to identify the exact class of the pixels. Cloud shadow detection is not the focus of this paper, therefore, this class is not labeled.

Fig. 8 shows the statistics of the Levir_CS dataset. The distribution of the label components is calculated from different views. As shown in Fig. 8 (a), in Levir CS, the background pixels occupy the most population (79.2%) while the snow occupies the least (2.2%). The cloud pixels occupy 18.6% of the whole pixel population. Fig. 8 (b)-(d) show the label components from the longitude, latitude and altitude views, respectively. From these figures, the following observations can be summarized:

- The pixel population of the three categories is very different in different locations.
- Clouds are common in different geolocations. For example, in North America, clouds appear in different kinds and forms (Sun et al., 2017).
- From the view of the longitude, it can be seen that the snow may be less likely to appear in the range of $-60^\circ \leqslant long \leqslant 30^\circ$ (Atlantic Ocean) (see Fig. 7(c,e) for examples), while in the area of $65^\circ \leqslant long \leqslant 100^\circ$, it is easier to find snow cover (see Fig. 7(d) for an example).
- From the view of the latitude, it can be seen that most of the snow appears in the high latitudes ($lat \geqslant 43^\circ$) (see Fig. 7(f) for an example). In the United States, the number of snow days is higher in the high latitude regions, according to (Tran et al., 2019). Besides, at high latitudes in the polar region, snow and ice does not melt in some seasons (Selkowitz and Forster, 2016). There is almost no snow covers the Equatorial regions ($-23.5^\circ \leqslant lat \leqslant 23.5^\circ$) (see Fig. 7(a,c,e) for examples).
- From the view of the altitude, it can be seen that the cloud percentage is higher in the area where the altitude is less than 500 meters (see Fig. 7(a,b,e) for examples) and the snow percentage gradually increases as the altitude increases (see Fig. 7(d,f) for examples). Usually, the high altitude area is mountainous area, and snow cover is regularly changed in seasons here (Wang et al., 2018). For the area with an altitude higher than 3400 meters, it is even easier to find snow than to find clouds (see Fig. 7(d) for an example).

From the above statistics, it can be seen that using geographic information in cloud and snow detection is of importance.

5. Experimental results and analysis

In this section, extensively evaluation are made on the proposed method and compare it with other state of the art ones. First the implementation details are introduced and how the experiments are set up are described. Then, the controlled experiments are conducted on multiple aspects of the method. Thirdly, qualitative and quantitative comparisons with other methods on Levir_CS are made. Finally, the transferability of the proposed GeoInfoNet by evaluating on other sensor data is tested on L8_Biome (Foga et al., 2017).

5.1. Experiment setup and implementation details

In this paper, all the deep learning models tested are implemented with PyTorch 1.0 on Ubuntu 16.04 with an NVIDIA Geforce GTX 1080Ti GPU card. The networks are trained by using the stochastic gradient descent method with an initial learning rate of 0.001. The learning rate decay policy is set to “poly” as (Wu and Shi, 2018) did, and the power parameter is set to 0.9. The number of iteration, the l_2 weight decay, and the momentum are set to 2×10^5 , 0.0001, and 0.9, respectively.

All backbone feature extraction models (including the Dense Feature Extraction module, VGG16 (Simonyan and Zisserman, 2015), ResNet101 (He et al., 2016), and DenseNet169 (Huang et al., 2017)) are pretrained on the Imagenet Dataset (Deng et al., 2009). The weights of the convolutional layers in the other components of the networks are initialized by the “msra” method (He et al., 2015). The number of feature maps provided by each additional convolution layer in the method is set to 64.

As informed in the above Section 4, the size of each image is 1320×1200 in both the training set and the test set. For the limitation of the GPU card memory, the training batch size is set to 4, and all the inputs are randomly cropped to 240×240 in the training phase. As for the testing phase, the input is cropped to 600×600 patches for evaluation and then these patches are combined together.

In the training phase, to increase the diversity of the images, data augmentation is performed by randomly rotating the inputs for $0^\circ, 90^\circ, 180^\circ, 270^\circ$. All the input bands are normalized to [0, 1]. In detail, for image bands, as the input images are 10-bit images, they are all divided by 1023. For the longitude bands, they are divided by 360 after added by 180. For the latitude bands, they are divided by 180 after added by 90. Finally, for the altitude bands, they are divided by 10000.

To evaluate the performance of different methods, three types of benchmark metrics are used in the experiments: *F1-Score* (*F1*), *Intersection-over-Union* (*IoU*), and *False Alarm Ratio* (*FAR*). These metrics are all widely used in cloud detection tasks (Chen et al., 2018; Wu and Shi, 2018; Zhan et al., 2017; Yan et al., 2018). In detail, for cloud or snow, *F1* is calculated as Eq. 7 shows.

$$F1 = \frac{2p \cdot r}{p + r}, \quad (7)$$

where *p* and *r* are calculated as Eq. 8 and Eq. 9,

$$p = \frac{N_{correct}}{N_{correct} + N_{false-alarm}}, \quad (8)$$

$$r = \frac{N_{correct}}{N_{ground-truth}}, \quad (9)$$

where $N_{correct}$ is the number of pixels of correct detection, $N_{false-alarm}$ is the number of pixels of false alarms and $N_{ground-truth}$ is the number of pixels of the certain type in the groundtruth images, respectively. Similarly, *IoU* of cloud or snow is calculated as Eq. 10 displays,

$$IoU = \frac{N_{correct}}{N_{ground-truth} + N_{false-alarm}}. \quad (10)$$

FAR is a type of benchmark to show the performance of all classes of detection (including cloud, snow and background), and it can be calculated as,

$$FAR = \frac{N_{wrong}}{N_{all}}, \quad (11)$$

where N_{wrong} is the number of pixels of wrong detection and N_{all} is the number of all pixels.

For the *F1* and the *IoU*, a higher score indicates better and the results on the cloud and the snow are recorded accordingly. In the following statistics, these two figures are differently recorded in different classes.

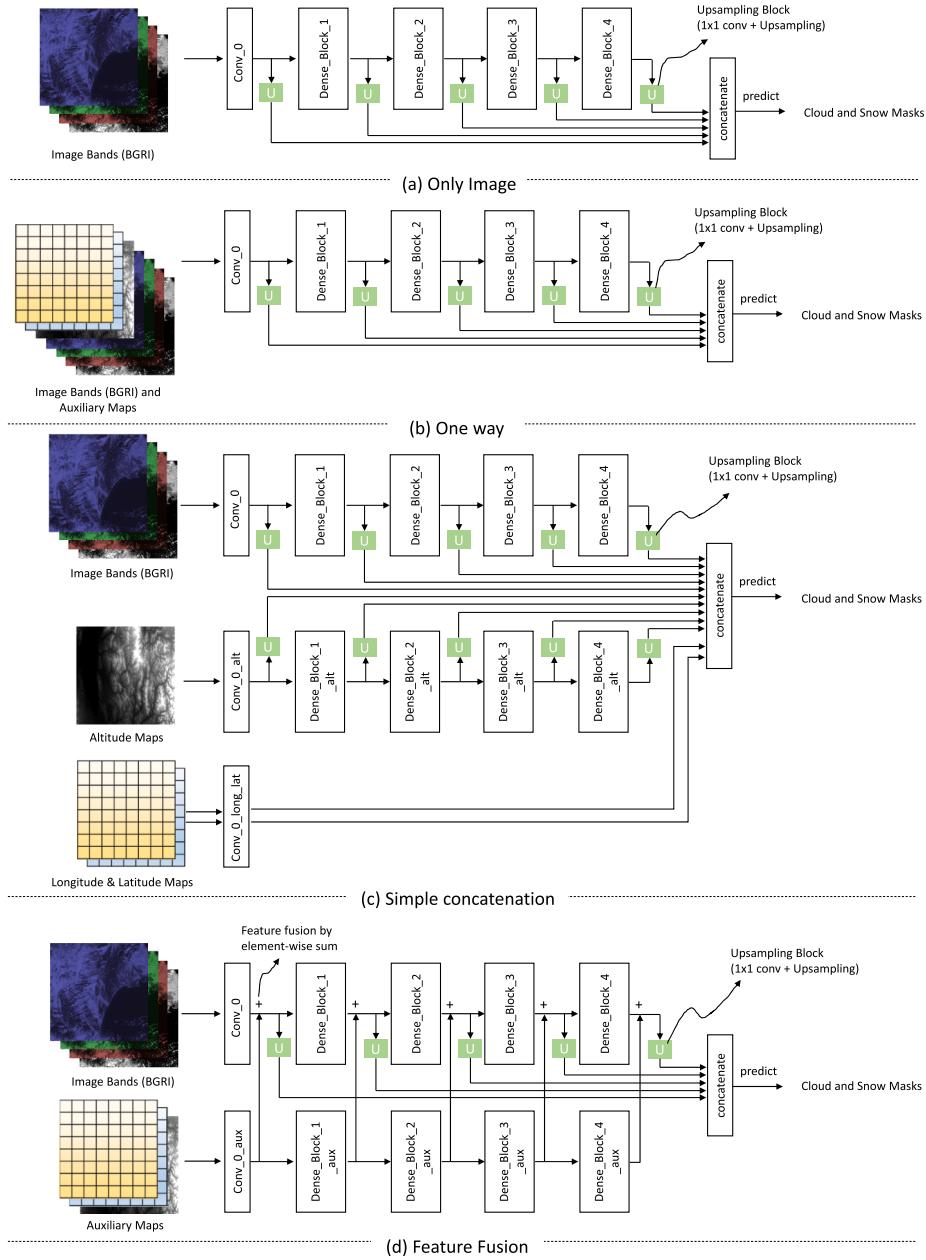


Fig. 9. Other possible network structures in the method in addition to the one used in Fig. 5: (a) Only Image. (b) One Way. (c) Simple Concatenation. (d) Feature Fusion (element-wise sum).

For the cloud, $F1_c$ and IoU_c are used to represent the capacity of cloud detection, while for the snow, $F1_s$ and IoU_s are used instead. For the FAR, a low score indicates better and the result on the whole test set of the Levir_CS is recorded.

Inference time is also evaluated per scene (image size: 1320×1200). In the testing phase, if the altitude map is used, it is directly loaded from the previous generated maps. This strategy is also used in the training phase. Therefore, it should be noted that the inference time does not include the time to create the corresponding altitude map. The time cost of making the altitude map has been introduced in the above Section 4 (45.62 s per scene) or other mentioned.

5.2. Controlled experiments

In this subsection, three types of controlled experiments are conducted, which focus on the verification of different technical components of the method: the network structure, auxiliary components

design, and the selection of backbone feature extractor. All these studies aim to find a reasonable design of the proposed method from different views.

5.2.1. Network structure design

It should be noticed that in addition to the network used in Fig. 5, there are also other possible structures can be chosen. These chooses are all suitable for the cloud and snow detection tasks but may have different accuracy performance. Fig. 9 presents four alternative choices on network configurations for the method.

Choice 1. In Fig. 9 (a), a “Only Image” structure is adopted in the cloud and snow detection tasks. In this structure, only image bands are processed. The cloud and snow masks are predicted by using the concatenated dense features from only the information of the image. This structure can be viewed as a baseline network structure since many strategies can be applied in this structure to improve the performance.

Choice 2. In Fig. 9 (b), a “One Way” structure is designed in a

Table 3

Quantitative results of four possible network configurations in the method. The result of F1, IoU, and FAR are recorded (%). The subscript “c” and “s” refer to the class of “cloud” and “snow”.

Network Structure	<i>F1_c</i>	<i>F1_s</i>	<i>IoU_c</i>	<i>IoU_s</i>	FAR	<i>T</i>
Simple Concatenation	Only Image	94.37	83.10	89.34	71.08	2.58
	One Way	94.62	83.94	89.79	72.32	2.44
	Fusion	94.71	85.74	89.95	75.03	2.40
	Dual Feature Concatenation	95.00	86.70	89.62	76.52	2.28
	Concatenation	95.15	87.80	90.75	78.26	2.20

straightforward way, which concatenates all the bands together before extracting dense features. The image information and the auxiliary information are thus processed together in only one network branch.

Choice 3. Another alternative design is “Simple Concatenation”, which is shown in Fig. 9 (c). The intuition behind this design is that the longitude and the latitude of the pixels do not need to be excessively

processed because their values within each scene will not change too much. Therefore, it is supposed that only the image-like bands need to be fed into the Dense Feature Extraction module, and the maps of the longitude and latitude are only one-layer convoluted before the feature concatenation.

Choice 4. The last choice is “Feature Fusion”, which is shown in Fig. 9 (d). In this choice, the same network structure is used as the default settings, while the only differences are: 1) the “feature concatenation” is changed to the “element-wise sum” operation, and 2) the features from the auxiliary branch within each stage are fused to the BGRI image branch before the further procedures.

Table 3 shows the comparison results of the four network structures in the method (the default structure and the *Choice 1–4*). Slight differences can be observed in the results of the different settings and the proposed *Dual Feature Concatenation* shown in Fig. 5 achieves the best in all metrics.

Table 4

Quantitative results of all possible combinations of auxiliary information in the cloud detection and snow detection tasks (%).

Longitude	Latitude	Altitude	<i>F1_c</i>	<i>F1_s</i>	<i>IoU_c</i>	<i>IoU_s</i>	FAR	<i>T</i>
			94.37	83.10	89.34	71.08	2.58	2.24s
✓			95.10	84.39	90.65	73.00	2.24	3.34s
	✓		94.85	84.88	90.20	73.73	2.37	3.28s
		✓	94.85	85.04	90.21	73.97	2.40	3.31s
✓	✓		95.05	85.08	90.57	74.03	2.27	3.34s
✓		✓	95.01	85.95	90.49	75.36	2.28	3.34s
	✓	✓	94.47	85.47	89.53	74.63	2.54	3.27s
✓	✓	✓	95.15	87.80	90.75	78.26	2.20	3.45s

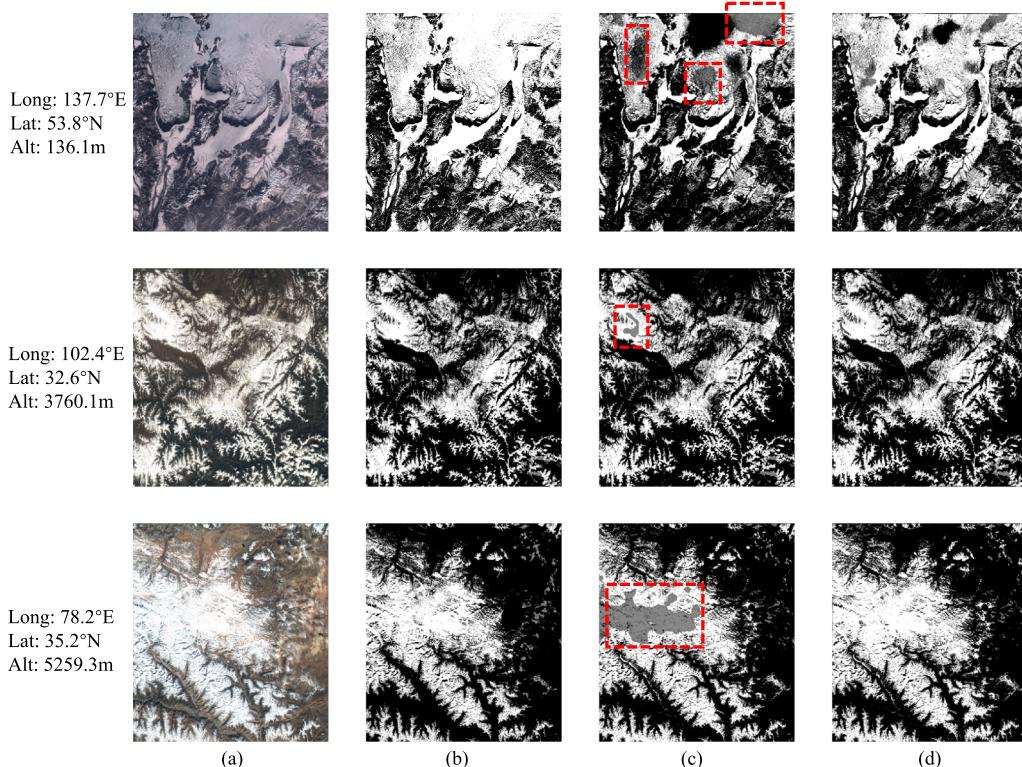


Fig. 10. Results of the method w/ and w/o using auxiliary information. (a) Input image. (b) Ground truth label. (c) Detection results w/o auxiliary information. (d) Detection results w/ auxiliary information. In (b)–(d), the white, grey, and black pixels represent the snow, cloud, and background, respectively. Red boxes show the false alarms of clouds which should be detected as snow, if auxiliary information is not used.

Table 5

Cloud and snow detection improvement on different locations. (%)

	w/o	w/	Δ	w/o	w/	Δ
	Cloud detection IoU_c			Snow detection IoU_s		
$long \leq -60^\circ$	92.19	93.93	1.74	75.76	83.91	8.15
$-60^\circ < long \leq 30^\circ$ ¹	88.13	89.43	1.30	-	-	-
$30^\circ < long \leq 65^\circ$	84.67	86.87	2.20	58.41	60.59	2.18
$65^\circ < long \leq 100^\circ$	83.43	85.54	2.11	75.18	79.53	4.35
$100^\circ < long \leq 120^\circ$	88.62	89.86	1.24	59.11	68.42	9.31
$long > 120^\circ$	90.74	91.08	0.34	70.85	79.17	8.32
$lat \leq -23.5^\circ$	86.29	87.59	1.30	76.96	86.54	9.58
$-23.5^\circ < lat \leq 23.5^\circ$ ²	92.07	92.89	0.82	-	-	-
$23.5^\circ < lat \leq 43^\circ$	89.04	90.57	1.53	75.58	78.52	2.94
$lat > 43^\circ$	84.44	86.98	2.54	69.51	77.96	8.45
$alt \leq 50 \text{ m}$	91.07	91.48	0.41	68.01	81.96	13.95
$50 \text{ m} < alt \leq 500 \text{ m}$	88.60	90.49	1.89	72.25	79.01	6.76
$500 \text{ m} < alt \leq 2100 \text{ m}$	88.19	89.88	1.69	62.48	72.43	9.95
$2100 \text{ m} < alt \leq 3400 \text{ m}$	90.46	93.23	2.77	72.79	81.04	8.25
$3400 \text{ m} < alt \leq 5500 \text{ m}$	88.84	90.27	1.43	82.08	83.22	1.14
$alt > 5500 \text{ m}$	68.47	84.60	16.13	70.70	77.44	6.74

¹ On the test set of Levir_CS, there is little snow in this longitude range, therefore, the snow detection results of this range are not able to be collected.

² On the test set of Levir_CS, there is little snow in this latitude range, therefore, the snow detection results of this range are not able to be collected.

5.2.2. Ablation studies on the auxiliary maps

In this experiment, ablation studies are conducted on three different types of geographic information: 1) longitude, 2) latitude, and 3) altitude. To show the importance of these auxiliary components, different combinations of them are evaluated. The network configurations for all the possible auxiliary information combinations are the same (described in the above Section 5.1), except for the number of input channel of the first convolution layer “Conv_0_aux”. The order of auxiliary information is always: altitude, longitude and latitude. For example, if only the information of latitude A_{lat} and altitude A_{alt} is used, the auxiliary map will be formed as $A = \text{concat}(A_{Alt}, A_{Lat})$. Table 4 shows the quantitative results of all the possible combinations of auxiliary information. The result indicates that the use of the auxiliary information is effective in the cloud and snow detection task. It can be seen that the integration of any of these auxiliary components brings noticeable improvement on the detection accuracy and the best result can be achieved when all the auxiliary components are integrated. In the cases that no altitude information is provided, adding the longitude map and the latitude map still can be beneficial for the cloud and snow detection.

Besides, detailed comparisons are also made on detection results about w/ and w/o auxiliary information. From Table 4, it can be noticed that with auxiliary information involved, the snow detection performance impressively improve. Fig. 10 is an illustration of w/ and w/o using auxiliary information. If auxiliary information is adopted, the snow area which is wrongly recognized as cloud will be obviously

reduced. It should be noted that snow detection is a very challenging task. As the first row of Fig. 10 illustrates, even with the auxiliary information, the snow area can be wrongly detected as cloud. This is why the performance of snow detection is lower than that of cloud. From Table 5, it can be seen that using auxiliary information can improve cloud and snow detection performance on different geographical location ranges. Specially, the performance of cloud detection raises if using the auxiliary information. For the region where the altitude is higher than 5500 m, the false alarms of clouds reduce as the last row of Fig. 10 implies. The performance of snow detection is generally improved in different locations. In the low altitude regions, the snow detection performance increases with a large margin, which can help to reduce the missing alarms of snow as the first row of Fig. 10 indicates. Therefore, the auxiliary geographic information helps the detection and does not hinder cloud or snow detection.

5.2.3. Evaluation on network backbone and loss function

For many deep learning based cloud and snow detection methods (Zhan et al., 2017; Yan et al., 2018; Wu and Shi, 2018; Shao et al., 2019; Yang et al., 2019; Li et al., 2019a), their backbone networks are built based on VGG (Simonyan and Zisserman, 2015) or ResNet (He et al., 2016). Therefore, in this experiment, different network backbones are evaluated for the task.

Besides, considering that the huge difference in the number of pixels between different classes, whether class-balancing will be beneficial to the detection accuracy is also evaluated. In this case, the weighted softmax loss is used, which is similar to (Xie and Tu, 2015), as the loss function of pixel, which is expressed as follows:

$$L_{\text{weighted}} = - \sum_{m=1}^3 \alpha_m y_m \log(P_m), \quad (12)$$

where α_m is the balancing weight for each class. α_m is set according to the pixel ratio of different classes, $\alpha_m = (n_m^{-1}) / (\sum_{i=1}^3 n_i^{-1})$, where n_m is the number of pixels belong to the class m in the dataset.

Table 6 shows the quantitative results of the method with different network backbones and losses. For the choice of the network backbones, the default design, i.e., the DenseNet169, obtains the best results. The accuracy increase is particularly significant for the snow detection task. As for the weighted loss, it can be observed that if the network backbone is VGG16 or ResNet101, the class-balancing can be useful. However, it does not help the detection if the network backbone is DenseNet169. As shown in Fig. 11, if the weighted loss is used, there will be more false alarms in the detection result, especially for the snow pixels. The above observations indicate that the backbone of DenseNet169 can be a robust network backbone and suffer little impact of imbalanced data. If class-balancing techniques are applied, the detection results may get worse. Therefore, the Dense Feature Extraction produces robust image features where the imbalanced data problems can be alleviated. This is the reason why class-balancing is not used in the default loss function.

Table 6

Quantitative results of the method with different network backbones and weighted loss (%).

Backbones	Weighted	$F1_c$	$F1_s$	IoU_c	IoU_s	FAR	T
VGG16	×	77.49	39.88	63.25	24.91	9.30	3.31s
VGG16	✓	88.48	57.10	79.33	39.95	5.20	3.31s
ResNet101	×	84.72	59.05	73.48	41.90	6.35	4.00s
ResNet101	✓	90.67	68.71	82.94	52.33	4.45	4.00s
DenseNet169	×	95.15	87.80	90.75	78.26	2.20	3.45s
DenseNet169	✓	93.98	85.18	88.65	74.20	2.81	3.45s

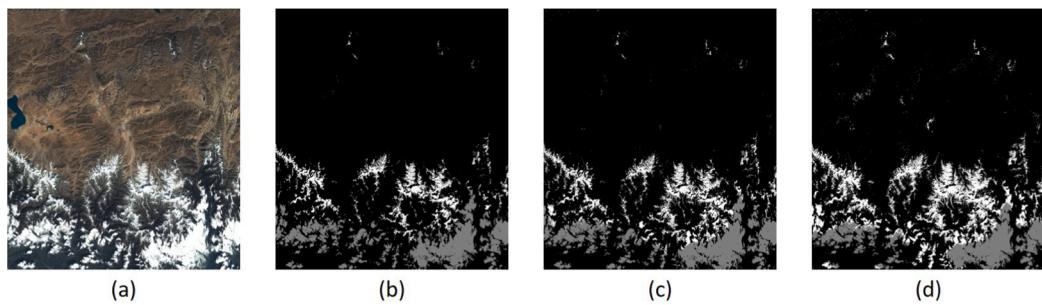


Fig. 11. Results of the method w/ and w/o using class-balancing in the loss function. (a) Input image. (b) Ground truth label. (c) Detection results w/o class-balancing. (d) Detection results w/ class-balancing. In (b)-(d), the white, grey, and black pixels represent the snow, cloud, and background, respectively. More false alarms are observed in the detection results when the class-balancing is used, especially for the snow pixels.

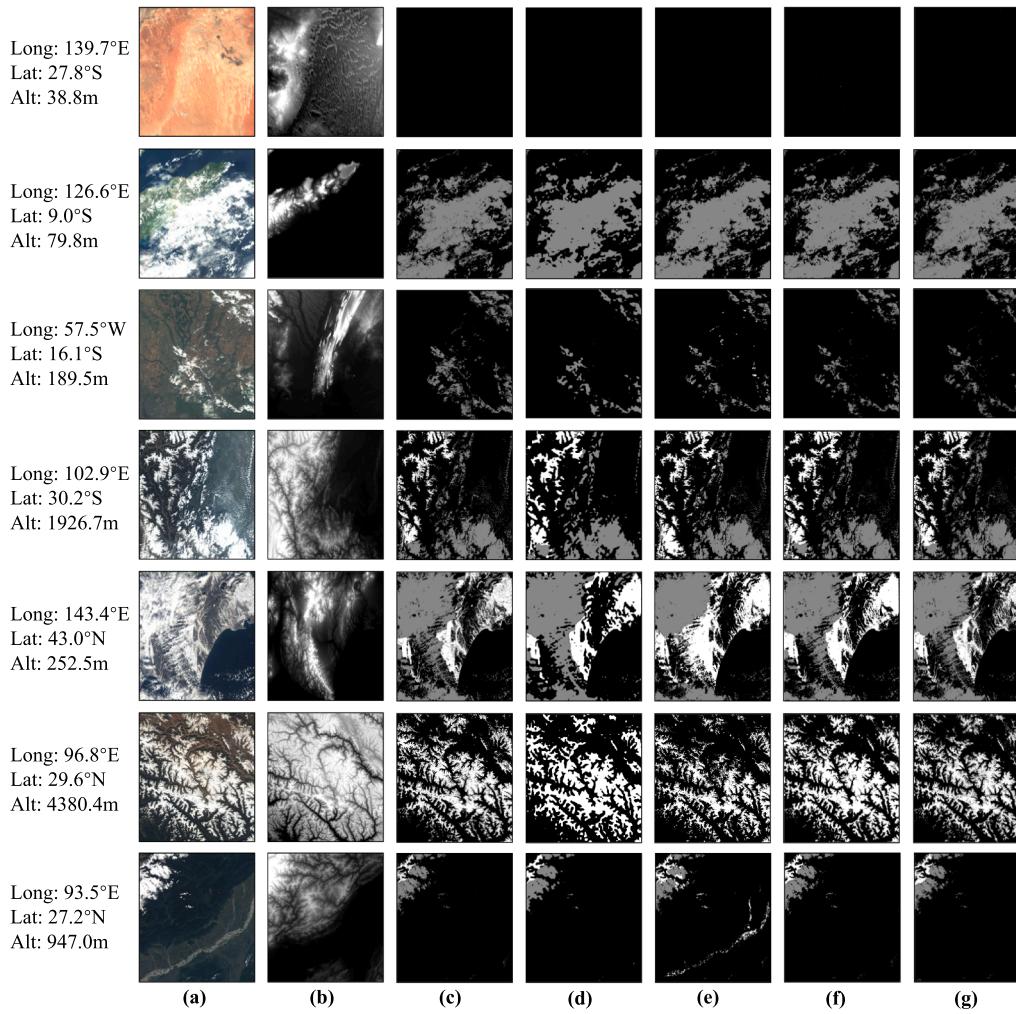


Fig. 12. Visual comparison of different cloud and snow detection methods. (a) Input image. (b) The corresponding altitude map. (c) Ground truth label. (d)-(g) The detection results of DCN (Zhan et al., 2017), DeepLabV3+ (Chen et al., 2018), GeoInfoNet4× (ours) and GeoInfoNet (ours), respectively. The white, grey, and black pixels represent the snow, cloud, and background. On the very left side of each row, the longitude and latitude of the central point and the mean altitude of the image is given.

5.3. Comparison with other methods

In this section, the method is compared with four state-of-the-art cloud and snow detection methods, the DCN (Zhan et al., 2017), FECN (Wu and Shi, 2018), cloudFCN (Francis et al., 2019) and cloudUNet (Jeppesen et al., 2019). The DCN (Zhan et al., 2017) is built based on the VGG16 backbone and produces the detection probability map from each level of the network. The final detection map of the DCN is obtained by summarizing all these probability maps. FECN (Wu and Shi, 2018) is also designed based on the VGG16 backbone. In this method, the features

from all stages are concatenated for the final prediction. The cloudFCN (Francis et al., 2019) is a cloud detection method based on fully convolutional neural networks. It modifies the original framework of FCN (Long et al., 2015) and can be successfully applied in cloud detection tasks. Similarly, cloudUNet (Jeppesen et al., 2019) is a method based on UNet (Ronneberger et al., 2015). In this framework, shallow features are reused in the upsampling procedures. Note that FECN, cloudFCN and cloudUNet are originally proposed for cloud detection, and in the experiment, we extend them for both cloud and snow detection by modifying the number of output classes from 2 (cloud and background)

Table 7

Quantitative comparisons of different methods on the Levir_CS dataset (%).

Methods	<i>F1_c</i>	<i>F1_s</i>	<i>IoU_c</i>	<i>IoU_s</i>	FAR	<i>T</i>
Scene Learning (An and Shi, 2015)	35.02	–	21.23	–	63.70	74.62 s
Coarse-to-Fine (Kang et al., 2018)	57.63	–	40.48	–	15.37	342 s
DCN (Zhan et al., 2017)	90.64	59.08	82.88	41.92	4.05	1.03 s
FECN (Wu and Shi, 2018)	89.72	60.25	81.36	43.12	4.33	2.17 s
cloudFCN (Francis et al., 2019)	93.01	75.50	86.94	60.64	3.06	2.59 s
cloudUNet (Jeppesen et al., 2019)	94.03	82.54	88.73	70.27	2.65	2.00 s
FCN-16s (Long et al., 2015)	79.19	69.32	65.54	53.04	8.62	0.86 s
DeepLabV3+ (Chen et al., 2018)	89.15	76.58	80.43	62.05	4.95	0.93 s
Only Image (ours)	94.37	83.10	89.34	71.08	2.58	2.24 s
GeoInfoNet4× (ours)	89.39	78.53	80.82	64.66	4.66	3.34 s
GeoInfoNet (ours)	95.15	87.80	90.74	78.26	2.20	3.45 s

to 3 (cloud, snow, and background).

Since the cloud and snow detection is essentially a semantic segmentation problem, the FCN-16s (Long et al., 2015) and DeepLabV3+ (Chen et al., 2018), which are two well-known image segmentation methods in computer vision, are also compared. FCN-16s (Long et al., 2015) uses the 16× downsampled feature maps for segmentation. DeepLabV3+ (Chen et al., 2018) integrates both low-level features (4× downsampled) and high-level features (produced by spatial pyramid pooling at the end of the network encoder and are 16× downsampled), and performs 4× upsampling on the prediction result. For a fair comparison, the network backbones used in these two methods are both DenseNet169, which is the same as used in the method. Besides, a low-resolution version of the method: “GeoInfoNet4×” is also experimented, which only uses 4× and larger times downsampled features for segmentation, for the comparison of the feature scales. In addition, traditional statistic methods based on machine learning Scene Learning (An and Shi, 2015) and Coarse-to-Fine (Kang et al., 2018) are also evaluated. These methods are used for cloud detection and cannot finish the snow detection tasks.

Visual comparison results of different detection methods are shown in Fig. 12. Table 7 shows their quantitative evaluation results. It can be seen that the GeoInfoNet method generates more accurate cloud and snow masks than other methods. Even without using higher resolution features, the GeoInfoNet4× still achieves satisfying performance (especially the snow detection). Compared with deep learning methods, traditional methods based on machine learning may not be proper in cloud detection tasks on global region data. Although the DCN (Zhan

Table 8

Quantitative evaluation results on the test set of L8_Biome (Foga et al., 2017) (%).

Networks	<i>F1_c</i>	<i>IoU_c</i>	FAR	<i>T</i>
Only Image (ours)	95.17	90.79	2.79	118.82 s
GeoInfoNet (ours)	95.84	92.01	2.43	188.00 s

et al., 2017) and the FECN (Wu and Shi, 2018) do not perform very well in snow detection (probably due to the imbalanced classes as we analyzed in Section 5.2), they still outperforms the FCN-16s (Long et al., 2015), the DeepLabV3+ (Chen et al., 2018), and even the GeoInfoNet4× in cloud detection. Based on the above observations, it can be concluded that for cloud and snow detection, utilizing all scales of features can obtain better detection results than those methods which only use lower resolution features.

5.4. Evaluation on other sensor data

In this subsection, the proposed GeoInfoNet is evaluated on other sensor data to test the transferability of the method. L8_Biome (Foga et al., 2017) is chosen in this part. As L8_Biome (Foga et al., 2017) does not contain snow information masks, only evaluate only the performance on cloud detection is evaluated. In L8_Biome (Foga et al., 2017), images are randomly divided into the training set and the testing set. The number of scenes of the training set is 77 while that of the testing set is 18. To generate the corresponding digital elevation map, the original images is first transferred in the World Geodetic System (WGS84) by using the GDAL library (Warmerdam, 2008). The mean area of the scenes in the test set is 55,126,569 pixels (around 7425 px × 7425 px). The digital elevation map generation costs 158.11s time per scene on average. Similar experiment settings to those introduced in Section 5.1 are adopted here, except for the number of iterations is set to 8×10^5 . Two network structures are evaluated, 1) the network with only image branch and 2) the proposed GeoInfoNet.

Fig. 13 illustrates the visual comparison results of different detection methods. Table 8 shows the quantitative evaluation results. It can be seen that the proposed method GeoInfoNet can be well applied to Landsat8. Besides, these results also prove the effectiveness of adding geographic information for detecting cloud.

6. Discussion

6.1. What prior information does our method learn?

In this section, an interesting question is raised: what kind of prior information does GeoInfoNet learn? A deep investigation has been made based on the method of Activation Maximization (Erhan et al., 2009).

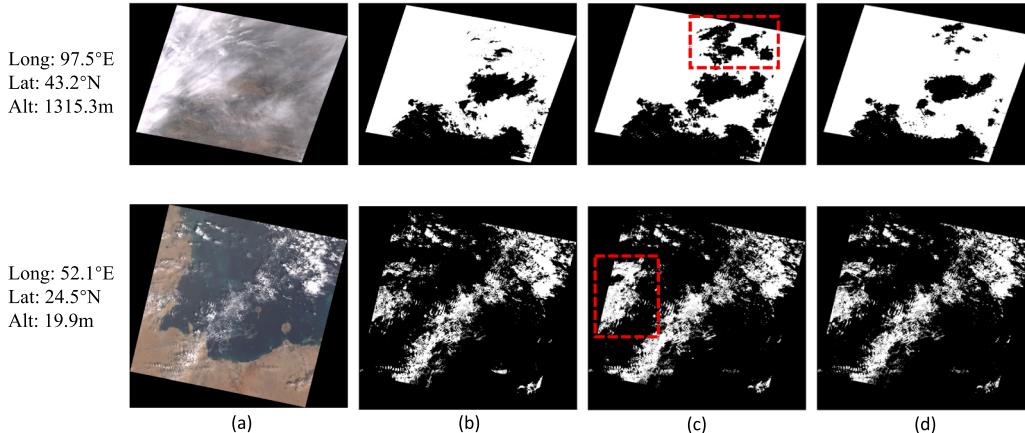


Fig. 13. Visual evaluation results on L8_Biome (Foga et al., 2017). (a) Input image. (b) Ground truth label. (c) The detection result of Only Image. (d) The detection result of GeoInfoNet. The white and black pixels represent the cloud (both thick and thin cloud) and the other categories (including the clear region, cloud shadow and the filling region). Red boxes shows the wrong-predictions (first row: more missing area; second row: more false alarms), if auxiliary information is not used. On the very left side of each row, the longitude and latitude of the central point and the mean altitude of the image is given.

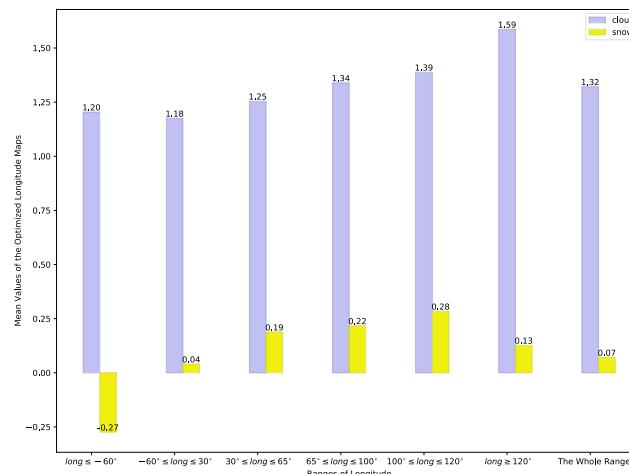


Fig. 14. The histogram of activation maximization results of the longitude maps.

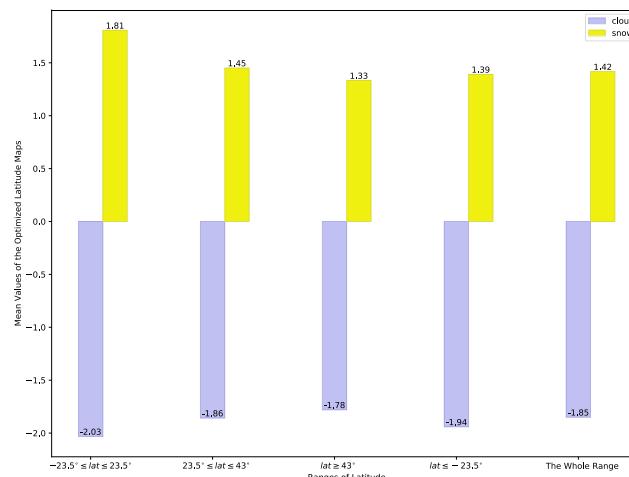


Fig. 15. The histogram of activation maximization results of the latitude maps.

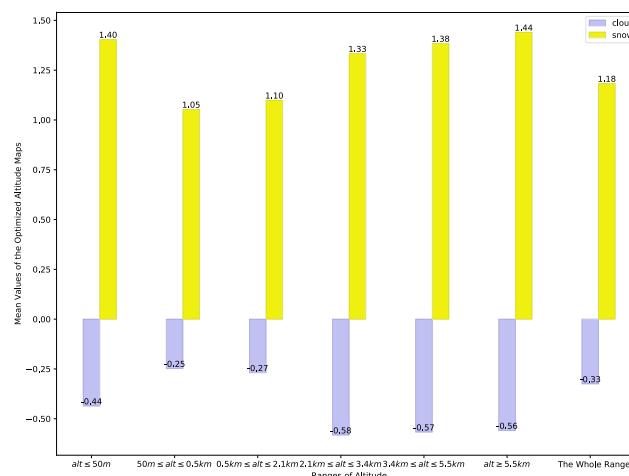


Fig. 16. The histogram of activation maximization results of the altitude maps.

This method was originally proposed to visualize the learned convolutional filters of the network by optimizing the feature maps. For a typical input U and a fixed network parameter Θ , the input U can be optimized as follows,

$$U^* = \operatorname{argmax}_{f_{ij}(\cdot, \cdot)} f_{ij}(\Theta, U), \text{s.t. } \|U\| \leq \rho, \quad (13)$$

where $f_{ij}(\cdot, \cdot)$ is an activation function of the input U and the network parameter Θ , given a convolution filter i from a given layer j in the network. Here, the parameters of GeoInfoNet are fixed, and the inputs of auxiliary branch A are optimized. For the details of optimization, Adam optimizer (Kingma and Ba, 2017) is used with the learning rate set to 0.1 and the number of iterations is set to 90. The activation function $f_{ij}(\cdot, \cdot)$ is set as the second or third channel of the outputs of the final score layer, which represent the score of the cloud class and snow class respectively. Through activation maximization, the inputs of auxiliary branch A will be changed through different activations, which represents the prior auxiliary information learned by GeoInfoNet.

Therefore, all the images in the testing sets can be optimized and the mean values of the optimized auxiliary information map have been calculated. Fig. 14, Fig. 15 and Fig. 16 illustrate the histograms of the mean values of the optimized auxiliary information maps in different location ranges. As the optimized values of the auxiliary maps have exceeded the real range values of geographic ranges, therefore, the activation maximization can only reflect the tendencies. From these figures, it can be seen that the method learns the prior knowledge that the cloud tends to appear in low altitude, low latitude and high longitude, while the trend for the snow is the opposite. This is consistent with the common sense because it is indeed easier to see snow in high altitudes. Besides, to some degree, it also grasps the data distribution shown in Fig. 8.

6.2. Analysis of the feature importance

In the method, since the features from all levels from both the image branch and the auxiliary branch are used, it is necessary to analyze how much these features of different levels and different branches contribute to the cloud and snow detection task.

In the proposed GeoInfoNet, the cloud and snow detection results are obtained according to the concatenated features M formed by different levels of features $M_{\text{img},0}, \dots, M_{\text{img},4}, M_{\text{aux},0}, \dots, M_{\text{aux},4}$ in both two branches (shown in Eq. 4). Therefore, the gradient of the prediction score S_t of a specific class t (cloud or snow in the topic) is computed with respect to M , and then multiply this gradient on the feature map to produce the “importance map” G of each pixel location on this feature map. The “importance” can be expressed as follows:

$$G_t = \sum_{i=1}^C \left(\frac{\partial S_t}{\partial M} M \right)^{(i)}, \quad (14)$$

where i is the channel index and C is the number of channels in the feature map M . Here, the gradient information conveys the neuron importance weight, therefore, by multiplying the gradient and the feature maps, the generated map G_t is a hot map which can show the most sensitive region of the specific class t grasped by the network. That is why G_t is named “importance map”. This process is similar to a well-known feature visualization method called Gradient Class Activation Map (Grad-CAM) (Selvaraju et al., 2017). The original Grad-CAM (Selvaraju et al., 2017) uses ReLU operation for it only interests in pixels with positive response to the specific class t . With ReLU function operated, those pixels with negative response are filtered. This operation is removed here because the negative values may reflect the negative tendency to this class (t), which may help to reduce the false alarms to this class.

Since the Eq. 14 is linear, the importance of different channel splits can also be computed from the equation. For example, the importance of the features from the image branch and the auxiliary branch can be efficiently computed by accumulating the above scores over the corresponding feature channels. Suppose G_t^{img} and G_t^{aux} represent for the feature importance of the two branches for the class t , and thus $G_t =$

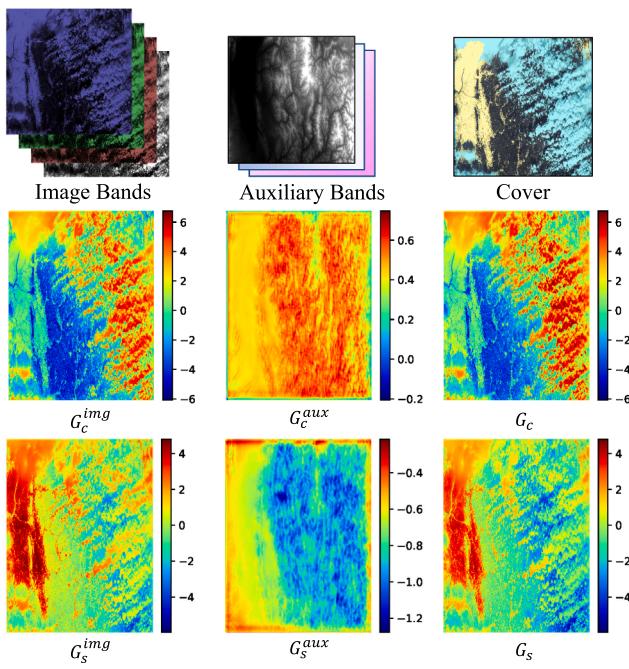


Fig. 17. Visualization results of the feature importance from the two different branches by using Eq. 14. The feature importance of cloud and snow are shown in the 2nd row and 3rd row, respectively. The image on the top-right is a cover image that represents the detection result, where light-blue represent cloud pixels while light-yellow represents snow pixels.

$G_t^{\text{img}} + G_t^{\text{aux}}$. Fig. 17 illustrates the importance maps target on the class of cloud and snow, which indicates that features from both two branches take effect in the computation of obtaining cloud and snow masks. Besides, the auxiliary information plays on an “auxiliary” role since the absolute values of G_c^{aux} and G_s^{aux} is much smaller than G_c^{img} and G_s^{img} . Therefore, the proposed GeoInfoNet still mainly relies on the image information and the auxiliary information helps the network improve the cloud and snow detection performance.

Besides, by using this method, the feature importance of different feature levels can also be easily obtained. Therefore, all the images in the

test set of Levir_CS are scanned and the importance of the different feature blocks in both branches is calculated. For each group of the features, an average importance score is computed, which is shown in Fig. 18. From this figure, it can be seen that in most of the feature groups (except for “Dense1” and “Dense2”), the importance of the image branch larger than that of the auxiliary branch in both cloud and snow detection tasks. This observation indicates that the image branch dominates the detection but the geographic information still contributes to the results to some degree. It can also be seen that in both branches, the importance of the very first convolution layer “Conv_0” and the very last convolution group “Dense_Block_4” are both very high, which shows that low-level features with high resolution and high-level features with low resolution are both crucial for cloud and snow detection tasks. As a comparison, for the groups from the “Dense_Block_1” to “Dense_Block_3”, the feature importance of both branches and both tasks is very close. Therefore, in the middle levels, features of both branches contribute to almost the same degree to the results.

7. Conclusion and future works

In this paper, a novel cloud and snow detection method is proposed for remote sensing images named “Geographic Information-drive Neural Networks (GeoInfoNet)”. Different from previous methods that simply perform detection solely based on the image data, the method integrates both the image and geographic information (altitude, latitude, and longitude) for training and detection. A large dataset for cloud and snow detection is also built, which contains 4,168 scenes and the corresponding geographic information. Extensive experiments verified the effectiveness of integrating geographic information for the cloud and snow detection tasks. The method outperforms other state-of-the-art methods with a large margin. Besides, the visualization is also presented to show what the method learns and how much the different parts of the network contribute to the detection tasks.

The future works include four parts. The first part is the improvement of the computational efficiency of the network. The second part of work is the integration of other types of geographic information (e.g., sun altitude angle, imaging time, temperature, etc). Third, cloud shadow will be focused in the future works, and this source of information will be integrated in the dataset. Finally, cloud and snow detection in time series (or multi-temporal cloud and snow detection) at specific locations will be investigated.

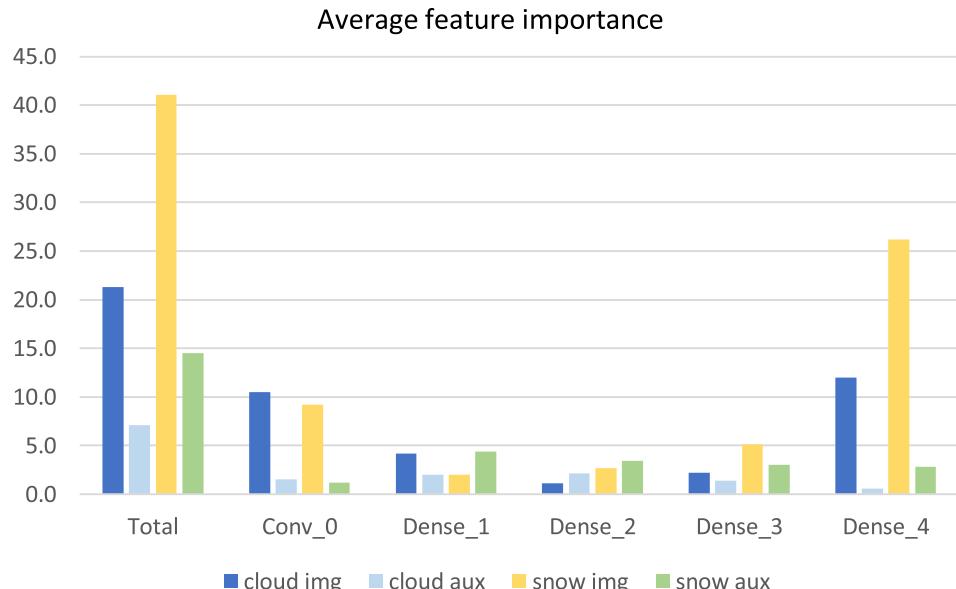


Fig. 18. The importance of the features from different levels and different branches for cloud detection and snow detection tasks. For each group of features, an average importance score on the test set of Levir_CS is computed.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

The work was supported by the National Key R&D Program of China under the Grant 2019YFC1510905, the National Natural Science Foundation of China under the Grant 61671037 and the Beijing Natural Science Foundation under the Grant 4192034.

References

- Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., Süsstrunk, S., 2012. Slic superpixels compared to state-of-the-art superpixel methods. *IEEE Trans Pattern Anal. Mach. Intell.* 34 (11), 2274–2282.
- An, Z., Shi, Z., 2015. Scene learning for cloud detection on remote-sensing images. *IEEE J. Sel. Top. Appl. Earth Observ. Remote Sens.* 8 (8), 4206–4222.
- Andersen, T., 1982. Operational snow mapping by satellites. In: Hydrological aspects of alpine and high mountain areas, Proceedings of the Exeter symposium, no. 138, pp. 149–154.
- Bian, J., Li, A., Jin, H., Zhao, W., Lei, G., Huang, C., 2014. Multi-temporal cloud and snow detection algorithm for the hj-1a/b ccd imagery of china. In: 2014 IEEE Geoscience and Remote Sensing Symposium. IEEE, pp. 501–504.
- Bian, J., Li, A., Liu, Q., Huang, C., 2016. Cloud and snow discrimination for ccd images of hj-1a/b constellation based on spectral signature and spatio-temporal context. *Remote Sens.* 8 (1), 31.
- Bi, J., Belle, J.H., Wang, Y., Lyapustin, A.I., Wildani, A., Liu, Y., 2019. Impacts of snow and cloud covers on satellite-derived pm2. 5 levels. *Remote Sens. Environ.* 221, 665–674.
- Breiman, L., 2001. Random forests. *Mach. Learn.* 45 (1), 5–32.
- Campbell, J.L., Mitchell, M.J., Groffman, P.M., Christenson, L.M., Hardy, J.P., 2005. Winter in northeastern north america: a critical period for ecological processes. *Front Ecol Environ.* 3 (6), 314–322.
- Chai, D., Newsam, S., Zhang, H.K., Qiu, Y., Huang, J., 2019. Cloud and cloud shadow detection in landsat imagery based on deep convolutional neural networks. *Remote Sens. Environ.* 225, 307–316.
- Chen, G., E., D., 2007. Support vector machines for cloud detection over ice-snow areas. *Geo. Spat. Inf. Sci.* 10 (2), 117–120.
- Chen, L.-C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H., 2018. Encoder-decoder with atrous separable convolution for semantic image segmentation. In: Proceedings of the European conference on computer vision (ECCV), pp. 801–818.
- Choi, H., Bindschadler, R., 2004. Cloud detection in landsat imagery of ice sheets using shadow matching technique and automatic normalized difference snow index threshold value decision. *Remote Sens. Environ.* 91 (2), 237–242.
- Cortes, C., Vapnik, V., 1995. Support-vector networks. *Mach. Learn.* 20 (3), 273–297.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L., 2009. Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. IEEE, pp. 248–255.
- Deng, C., Li, Z., Wang, W., Wang, S., Tang, L., Bovik, A.C., 2018. Cloud detection in satellite images based on natural scene statistics and gabor features. *IEEE Geosci. Remote Sens. Lett.* 16 (4), 608–612.
- Erhan, D., Bengio, Y., Courville, A., Vincent, P., 2009. Visualizing higher-layer features of a deep network. *Univ. Montreal 1341* (3), 1.
- Foga, S., Scaramuzza, P.L., Guo, S., Zhu, Z., Dilley Jr, R.D., Beckmann, T., Schmidt, G.L., Dwyer, J.L., Hughes, M.J., Laue, B., 2017. Cloud detection algorithm comparison and validation for operational landsat data products. *Remote Sens. Environ.* 194, 379–390.
- Francis, A., Sidiroopoulos, P., Muller, J.-P., 2019. Cloudfcn: Accurate and robust cloud detection for satellite imagery with deep learning. *Remote Sens.* 11 (19), 2312.
- Glorot, X., Bordes, A., Bengio, Y., 2011. Deep sparse rectifier neural networks. In: Proceedings of the fourteenth international conference on artificial intelligence and statistics, pp. 315–323.
- Hagolle, O., Huc, M., Pascual, D.V., Dedieu, G., 2010. A multi-temporal method for cloud detection, applied to formosat-2, venus, landsat and sentinel-2 images. *Remote Sens. Environ.* 114 (8), 1747–1755.
- He, K., Zhang, X., Ren, S., Sun, J., 2015. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: Proceedings of the IEEE international conference on computer vision, pp. 1026–1034.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778.
- Hollstein, A., Segl, K., Guanter, L., Brell, M., Enesco, M., 2016. Ready-to-use methods for the detection of clouds, cirrus, snow, shadow, water and clear sky pixels in sentinel-2 msi images. *Remote Sens.* 8 (8), 666.
- Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q., 2017. Densely connected convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4700–4708.
- Ioffe, S., Szegedy, C., 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift, arXiv preprint arXiv:1502.03167. Submission date: 2nd March, 2015.
- Irish, R.R., 2000. Landsat 7 automatic cloud cover assessment. In: Algorithms for Multispectral, Hyperspectral, and Ultraspectral Imagery VI, vol. 4049. International Society for Optics and Photonics, pp. 348–355.
- Irish, R.R., Barker, J.L., Goward, S.N., Arvidson, T., 2006. Characterization of the landsat-7 etm+ automated cloud-cover assessment (acca) algorithm. *Photogramm. Eng. Remote Sens.* 72 (10), 1179–1188.
- Jain, A.K., Farrokhnia, F., 1990. Unsupervised texture segmentation using gabor filters. In: 1990 IEEE international conference on systems, man, and cybernetics conference proceedings. IEEE, pp. 14–19.
- Jeppesen, J.H., Jacobsen, R.H., Inceoglu, F., Toftegaard, T.S., 2019. A cloud detection algorithm for satellite imagery based on deep learning. *Remote Sens. Environ.* 229, 247–259.
- Kang, X., Gao, G., Hao, Q., Li, S., 2018. A coarse-to-fine method for cloud detection in remote sensing images. *IEEE Geosci. Remote Sens. Lett.* 16 (1), 110–114.
- Kingma, D.P., Ba, J., 2017. Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980. Submission date: 30th Jan, 2017.
- Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems, pp. 1097–1105.
- Le Goff, M., Tournet, J.-Y., Wendt, H., Ortner, M., Spigai, M., 2017. Deep learning for cloud detection. In: ICPRs (8th International Conference of Pattern Recognition Systems), IET.
- Li, X., Shen, H., Zhang, L., Zhang, H., Yuan, Q., Yang, G., 2014. Recovering quantitative remote sensing products contaminated by thick clouds and shadows using multitemporal dictionary learning. *IEEE Trans. Geosci. Remote Sens.* 52 (11), 7086–7098.
- Li, P., Dong, L., Xiao, H., Xu, M., 2015. A cloud image detection method based on svm vector machine. *Neurocomputing* 169, 34–42.
- Li, Z., Shen, H., Li, H., Xia, G., Gamba, P., Zhang, L., 2017. Multi-feature combined cloud and cloud shadow detection in gaofen-1 wide field of view imagery. *Remote Sens. Environ.* 191, 342–358.
- Li, X., Jing, Y., Shen, H., Zhang, L., 2019b. The recent developments in cloud removal approaches of modis snow cover product. *Hydrocl. Earth Syst. Sci.* 23 (5).
- Li, Z., Shen, H., Cheng, Q., Liu, Y., You, S., He, Z., 2019a. Deep learning based cloud detection for medium and high resolution remote sensing images of different sensors. *ISPRS J. Photogramm. Remote Sens.* 150, 197–212.
- Li, W., Zou, Z., Shi, Z., 2020. Deep matting for cloud detection in remote sensing images. *IEEE Trans. Geosci. Remote Sens.* 58 (12), 8490–8502.
- Lin, H., Shi, Z., Zou, Z., 2017a. Maritime semantic labeling of optical remote sensing images with multi-scale fully convolutional network. *Remote Sens.* 9 (5), 480.
- Lin, H., Shi, Z., Zou, Z., 2017b. Fully convolutional network with task partitioning for inshore ship detection in optical remote sensing images. *IEEE Geosci. Remote Sens. Lett.* 14 (10), 1665–1669.
- Long, J., Shelhamer, E., Darrell, T., 2015. Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 3431–3440.
- Lu, J., Wang, Y., Zhu, Y., Ji, X., Xing, T., Li, W., Zomaya, A.Y., 2019. P.segnet and np-segnet: New neural network architectures for cloud recognition of remote sensing images. *IEEE Access* 7, 87323–87333.
- Mateo-García, G., Gómez-Chova, L., Camps-Valls, G., 2017. Convolutional neural networks for multispectral image cloud masking. In: 2017 IEEE International Geoscience and Remote Sensing Symposium (IGARSS). IEEE, pp. 2255–2258.
- Mehrotra, R., Namuduri, K.R., Ranganathan, N., 1992. Gabor filter-based edge detection. *Pattern Recognit.* 25 (12), 1479–1494.
- Munkres, J., 1957. Algorithms for the assignment and transportation problems. *J. Soc. Ind. Appl. Math.* 5 (1), 32–38.
- Musial, J.P., Hüslér, F., Stütterlin, M.B., Neuhaus, C., Wunderle, S., 2014. Probabilistic approach to cloud and snow detection on advanced very high resolution radiometer (avhr) imagery. *Atmos. Meas. Tech. (AMT)* 7 (3), 799–822.
- Osada, R., Funkhouser, T., Chazelle, B., Dobkin, D., 2002. Shape distributions. *ACM Trans. Gr.* 21 (4), 807–832.
- Qiu, S., He, B., Zhu, Z., Liao, Z., Quan, X., 2017. Improving fmask cloud and cloud shadow detection in mountainous area for landsats 4–8 images. *Remote Sens. Environ.* 199, 107–119.
- Qiu, S., Zhu, Z., He, B., 2019. Fmask 4.0: Improved cloud and cloud shadow detection in landsats 4–8 and sentinel-2 imagery. *Remote Sens. Environ.* 231, 111205.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. Springer, pp. 234–241.
- Scaramuzza, P.L., Bouchard, M.A., Dwyer, J.L., 2011. Development of the landsat data continuity mission cloud-cover assessment algorithms. *IEEE Trans. Geosci. Remote Sens.* 50 (4), 1140–1154.
- Selkowitz, D.J., Forster, R.R., 2016. An automated approach for mapping persistent ice and snow cover over high latitude regions. *Remote Sens.* 8 (1), 16.
- Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D., 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In: Proceedings of the IEEE international conference on computer vision, pp. 618–626.
- Shao, Z., Pan, Y., Diao, C., Cai, J., 2019. Cloud detection in remote sensing images based on multiscale features-convolutional neural network. *IEEE Trans. Geosci. Remote Sens.* 57 (6), 4062–4076.
- Shi, M., Xie, F., Zi, Y., Yin, J., 2016. Cloud detection of remote sensing images by deep learning. In: 2016 IEEE International Geoscience and Remote Sensing Symposium (IGARSS). IEEE, pp. 701–704.

- Shi, T., Xu, Q., Zou, Z., Shi, Z., 2018. Automatic raft labeling for remote sensing images via dual-scale homogeneous convolutional neural network. *IEEE Trans. Image Process.* 10 (7), 1130.
- Simonyan, K., Zisserman, A., 2015. Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556. Submission date: 10th April, 2015.
- Srivastava, A.N., Stroeve, J., 2003. Onboard detection of snow, ice, clouds and other geophysical processes using kernel methods. In: Proceedings of the ICML, Citeseer, vol. 3.
- Sun, L., Mi, X., Wei, J., Wang, J., Tian, X., Yu, H., Gan, P., 2017. A cloud detection algorithm-generating method for remote sensing data at visible to short-wave infrared wavelengths. *ISPRS J. Photogramm. Remote Sens.* 124, 70–88.
- Tran, H., Nguyen, P., Ombadi, M., Hsu, K.-L., Sorooshian, S., Qing, X., 2019. A cloud-free modis snow cover dataset for the contiguous united states from 2000 to 2017. *Sci. Data* 6, 180300.
- Wang, X., Wang, J., Che, T., Huang, X., Hao, X., Li, H., 2018. Snow cover mapping for complex mountainous forested environments based on a multi-index technique. *IEEE J. Sel. Top. Appl. Earth Observ. Remote Sens.* 11 (5), 1433–1441.
- Warmerdam, F., 2008. The geospatial data abstraction library. In: Open source approaches in spatial data handling. Springer, pp. 87–104.
- Weldon, T.P., Higgins, W.E., Dunn, D.F., 1996. Efficient gabor filter design for texture segmentation. *Pattern Recognit.* 29 (12), 2005–2015.
- Wieland, M., Li, Y., Martinis, S., 2019. Multi-sensor cloud and cloud shadow segmentation with a convolutional neural network. *Remote Sens. Environ.* 230, 111203.
- Wu, X., Shi, Z., 2018. Utilizing multilevel features for cloud detection on satellite imagery. *Remote Sens.* 10 (11), 1853.
- Xie, S., Tu, Z., 2015. Holistically-nested edge detection. In: Proceedings of the IEEE international conference on computer vision, pp. 1395–1403.
- Xie, F., Shi, M., Shi, Z., Yin, J., Zhao, D., 2017. Multilevel cloud detection in remote sensing images based on deep learning. *IEEE J. Sel. Top. Appl. Earth Observ. Remote Sens.* 10 (8), 3631–3640.
- Yan, Z., Yan, M., Sun, H., Fu, K., Hong, J., Sun, J., Zhang, Y., Sun, X., 2018. Cloud and cloud shadow detection using multilevel feature fused segmentation network. *IEEE Geosci. Remote Sens. Lett.* 15 (10), 1600–1604.
- Yang, J., Guo, J., Yue, H., Liu, Z., Hu, H., Li, K., 2019. Cdnet: Cnn-based cloud detection for remote sensing imagery. *IEEE Trans. Geosci. Remote Sens.* 57 (8), 6195–6211.
- Zhan, Y., Wang, J., Shi, J., Cheng, G., Yao, L., Sun, W., 2017. Distinguishing cloud and snow in satellite images via deep convolutional network. *IEEE Geosci. Remote Sens. Lett.* 14 (10), 1785–1789.
- Zhang, Q., Xiao, C., 2014. Cloud detection of rgb color aerial photographs by progressive refinement scheme. *IEEE Trans. Geosci. Remote Sens.* 52 (11), 7264–7275.
- Zhao, S., Yu, T., Meng, Q., Zhou, Q., Wang, F., Wang, L., Hu, Y., 2010. Gdal-based extend arcgis engine's support for hdf file format. In: 2010 18th International Conference on Geoinformatics. IEEE, pp. 1–3.
- Zhong, B., Chen, W., Wu, S., Hu, L., Luo, X., Liu, Q., 2017. A cloud detection method based on relationship between objects of cloud and cloud-shadow for chinese moderate to high resolution satellite imagery. *IEEE J. Sel. Top. Appl. Earth Observ. Remote Sens.* 10 (11), 4898–4908.
- Zhu, X., Helmer, E.H., 2018. An automatic method for screening clouds and cloud shadows in optical satellite image time series in cloudy regions. *Remote Sens. Environ.* 214, 135–153.
- Zhu, Z., Woodcock, C.E., 2012. Object-based cloud and cloud shadow detection in landsat imagery. *Remote Sens. Environ.* 118, 83–94.
- Zhu, Z., Woodcock, C.E., 2014. Automated cloud, cloud shadow, and snow detection in multitemporal landsat data: An algorithm designed specifically for monitoring land cover change. *Remote Sens. Environ.* 152, 217–234.
- Zi, Y., Xie, F., Jiang, Z., 2018. A cloud detection method for landsat 8 images based on peanet. *Remote Sens.* 10 (6), 877.
- Zou, Z., Shi, Z., 2016. Ship detection in spaceborne optical image with svd networks. *IEEE Trans. Geosci. Remote Sens.* 54 (10), 5832–5845.
- Zou, Z., Shi, Z., 2017. Random access memories: A new paradigm for target detection in high resolution aerial remote sensing images. *IEEE Trans. Image Process.* 27 (3), 1100–1111.
- Zou, Z., Li, W., Shi, T., Shi, Z., Ye, J., 2019. Generative adversarial training for weakly supervised cloud matting. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 201–210.