

1 Improved RepVGG Ground-Based Cloud Image Classification 2 with Attention Convolution

3 Chaojun Shi^{1, 2}, Leile Han¹, Ke Zhang^{1, 2}, Hongyin Xiang^{1, 2}, Xingkuan Li¹, Zibo Su¹, and Xian
4 Zheng¹

5 ¹Department of Electronic and Communication Engineering, North China Electric Power University, Baoding 071003,
6 China

7 ²Hebei Key Laboratory of Power Internet of Things Technology, North China Electric Power University, Baoding
8 Hebei 071003, China

9 *Correspondence to:* Chaojun Shi (scj@ncepu.edu.cn) and Hongyin Xiang (66283880@qq.com)

10 **Abstract.** Atmospheric clouds greatly impact the Earth's radiation, hydrological cycle, and climate change. Accurate
11 automatic recognition of cloud shape based on ground-based cloud image is helpful to analyze solar irradiance, water
12 vapor content, and atmospheric motion, and then predict photovoltaic power, weather trends, and severe weather
13 changes. However, the appearance of clouds is changeable and diverse, and its classification is still challenging. In
14 recent years, convolution neural network (CNN) has made great achievements in ground-based cloud image
15 classification. However, traditional CNNs poorly associate long-distance clouds, making the extraction of global
16 features of cloud images quite problematic. This study attempts to mitigate this problem by elaborating a ground-
17 based cloud image classification method based on the improved RepVGG convolution neural network and attention
18 mechanism. Firstly, the proposed method increases the RepVGG residual branch and obtains more local detail features
19 of cloud images through small convolution kernels. Secondly, an improved channel attention module is embedded
20 after the residual branch fusion, effectively extracting the global features of cloud images. Finally, the linear classifier
21 is used to classify the ground cloud images. Finally, the warm-up method is applied to optimize the learning rate in
22 the training stage of the proposed method, making it lightweight in the inference stage and thus avoiding overfitting
23 and accelerating the model's convergence. The proposed method is validated on MGCD and GRSCD ground-based
24 cloud image datasets containing 7 cloud categories, with the respective classification accuracy rate values of 98.15%
25 and 98.07%, outperforming those of ten most advanced methods used as the reference. The results obtained are
26 considered instrumental in ground-based cloud image classification.

27 1. Introduction

28 In meteorology, cloud is an aerosol consisting of a visible mass of water droplets, ice crystals, their aggregates or
29 other particles suspended in the atmosphere. Clouds of different types cover over 70% of the Earth surface (Qu et al.,
30 2021; Gyasi and Swarnalatha, 2023; Fabel et al., 2022). Cloud analysis plays a crucial role in meteorological

31 observation because clouds can affect the Earth's water cycle, climate change, and solar irradiance (Gorodetskaya et
32 al., 2015; Goren et al., 2018; Zheng et al., 2019). Cloud observation methods mainly include satellite observation
33 (Norris et al., 2016; Zhong et al., 2017; Li et al., 2023) and ground observation (Calbó and Sabburg, 2008; Nouri et
34 al., 2019; Lin et al., 2023). Satellite observation refers to the distribution, movement, and change of clouds observed
35 by high-resolution remote sensing satellites from above. When observing local sky regions, satellite observations have
36 low performance and are unable to obtain sufficient resolution to describe the characteristics of different cloud layers
37 in detail (Long et al., 2023; Sarukkai et al., 2020). Compared with satellite observation, ground-based observation
38 opens up a new way to monitor and understand regional sky conditions. Typical ground-based cloud observation
39 instruments include All-Sky Imager (ASI) (Shi et al., 2019; Cazorla et al., 2008), Total Sky Imager (TSI) (Long et al.,
40 2006; Tang et al., 2021), etc. The relevant equipment and ground-based cloud images are shown in Figure 1.



41

42 **Figure 1: Two kinds of ground-based cloud images and their observation equipment: (a) ASI ground-based cloud image**
43 **and its observation equipment (Cazorla et al., 2008; Shi et al., 2019); (b) TSI ground-based cloud image and its observation**
44 **equipment (Long et al., 2006).**

45 Ground-based cloud observation can obtain more obvious cloud characteristics by observing the information at the
46 bottom of the cloud, which is conducive to assisting the prediction of local photovoltaic power generation. Clouds
47 play an important role in maintaining the atmospheric radiation balance by absorbing short-wave and the ground not to
48 solar radiation (Taravat et al., 2015). Pv power prediction is affected by multiple factors such as cloud genus, cloud
49 cover change, solar irradiance, and solar cell performance in local areas, among which cloud genus is an important
50 factor affecting PV power prediction (Zhu et al., 2022). Therefore, it is of great significance to accurately obtain sky
51 cloud information through cloud observation and then accurately classify clouds for accurate prediction of
52 photovoltaic power generation (Alonso-Montesinos et al., 2016). The traditional ground-based cloud observation
53 method is mainly visual observation, which relies heavily on the experience of observers, cannot achieve
54 standardization. Therefore, ground-based cloud automatic observation has been widely concerned by scholars. In
55 recent years, with the development of digital image acquisition devices, many ground-based whole-sky cloud image
56 acquisition devices have emerged the world, providing massive data support for automatic ground-based cloud
57 observation (Pfister et al., 2003).

58 Ground-based cloud image classification is an important part of the foundation of automatic cloud observation and
59 is the key to climate change and photovoltaic power prediction. The classification of ground-based cloud images
60 mainly classifies each cloud image taken from the ground into the corresponding cloud genus by extracting cloud
61 image features, such as cirrus, cumulus, stratus, nimbostratus, etc. According to different cloud image feature
62 extraction methods, the ground-based cloud image classification method is divided into based on traditional machine
63 learning method and based on deep learning method (Simonyan and Zisserman, 2015; Krizhevsky et al., 2017; Hu et
64 al., 2018). Most of the ground-based cloud image classification methods based on traditional machine learning classify
65 cloud images by artificially designing cloud image features, while the ground-based cloud image classification
66 methods based on deep learning mainly classify cloud images through self-learning cloud image features of deep
67 neural network (DNN) (Wu et al., 2019).

68 Early ground-based cloud image classification studies relied on manual classification methods, which focused on
69 features such as texture, structure, and color, combined with traditional machine learning methods to classify ground-
70 based cloud images. These methods include a decision tree, K-nearest neighbor (KNN) classifier, support vector
71 machine (SVM), etc. (Singh and Glennen, 2005) proposed a method for automatically training the texture function of
72 a cloud classifier. In this method, five feature extraction methods including autocorrelation, co-occurrence matrix,
73 edge frequency, Laws texture analysis, and original length are used respectively. Compared with other cloud
74 classification methods, this method has the advantages of high accuracy and fast classification speed, but its
75 classification ability for mixed clouds is insufficient. (Heinle et al., 2010) described cloud images by using spectral
76 features (mean value, standard deviation, skewness, and difference) and texture features (energy, entropy, contrast,
77 homogeneity, and cloud cover), and combined with a KNN classifier, divided ground cloud images into seven
78 categories. In addition, (Zhuo et al., 2014) reported that the spatial distribution of contour lines could represent the
79 structural information of cloud shapes, used the central description pyramid to simultaneously extract the texture and
80 structural features of ground-based cloud images, and used SVM and KNN to classify cloud images. It can be seen
81 that the traditional classification method of ground-based cloud images based on machine learning mainly uses hand-
82 designed texture, structure, color, shape, and other features to extract, and obtains high-dimensional feature expression
83 of ground-based cloud images through single feature or fusion feature. Traditional machine learning methods mostly
84 describe the features from the perspective of digital signal analysis and mathematical statistics, but ignore the
85 representation and interpretation of the visual features of the cloud image itself.

86 In recent years, under the background of cross-integration of different disciplines and artificial intelligence, the
87 ground-based cloud image classification method based on deep learning has become a research hotspot with its
88 superior classification performance. Aiming at the unique characteristics of ground-based cloud images, (Shi et al.,
89 2017) proposed Deep Convolutional Activations-Based Features (DCAFs) to classify ground-based cloud images,
90 and the results are better than the artificially designed cloud image features. Alternatively, (Ye et al., 2017) used CNN
91 to extract cloud image features and proposed a local pattern mining method based on ground-based cloud images to

92 optimize the local features of cloud images and improve the classification accuracy of cloud images. (Zhang et al.,
93 2018a) put the wake cloud as a new genus of cloud into the ground-based cloud image database for the first time,
94 proposed a simple convolutional neural network model called CloudNet, and applied it to the ground-based cloud
95 image classification task, effectively improving the accuracy of ground-based cloud image classification. More
96 recently, (Wang et al., 2020) proposed the CloudA network, an optimized iteration of the AlexNet convolutional
97 neural network, which reduces the number of parameters through a simplified network architecture. The classification
98 accuracy on the Singapore Whole-Sky Imaging Categories (SWIMCAT) ground-based cloud image dataset exceeded
99 the traditional ground-based cloud image classification methods. (Liu et al., 2020b) proposed Multi-Evidence and
100 Multi-Modal Fusion Networks (MMFN) by fusing heterogeneous features, local visual features, and multi-mode
101 information, which significantly improves the classification accuracy of cloud images. Aiming at the problem that the
102 traditional neural network has insufficient ability to classify the ground-based cloud images within and between genera,
103 (Zhu et al., 2022) proposed to use of an improved combined convolutional neural network to classify the cloud images,
104 and the classification accuracy is greatly improved compared with the traditional neural network. Alternatively, (Yu
105 et al., 2021) used two sub-convolutional neural networks to extract features of ground-based cloud images and used
106 weighted sparse representation coding to classify them, which solved the problem of occlusion in multi-mode ground-
107 based cloud image data and greatly improved the robustness of cloud images classification. (Liu et al., 2020a)
108 introduced a ground-based cloud image classification method based on a graph convolution network (GCN). However,
109 the weight assigned by GCN failed to accurately reflect the importance of connection nodes, thus reducing the
110 discrimination of aggregated cloud image features. To make up for this deficiency, (Liu et al., 2022) proposed a context
111 attention network for ground-based cloud classification and publicly released a new cloud classification dataset. In
112 addition, (Liu et al., 2020c) further combined CNN and GCN to propose a multimodal ground-based cloud image
113 classification method based on heterogeneous deep feature learning. Alternatively, (Wang et al., 2021) elaborated a
114 ground-based cloud image classification method based on Transfer Convolutional Neural Network (TCNN) by
115 combining deep learning and transfer learning. (Li et al., 2022) further enhanced the classification performance of
116 ground-based cloud images based on the improved Vision Transformer combined with the EfficientNet-CNN. The
117 performance of the above-mentioned ground-based cloud image classification methods based on deep learning has
118 significantly improved compared to traditional machine learning methods.

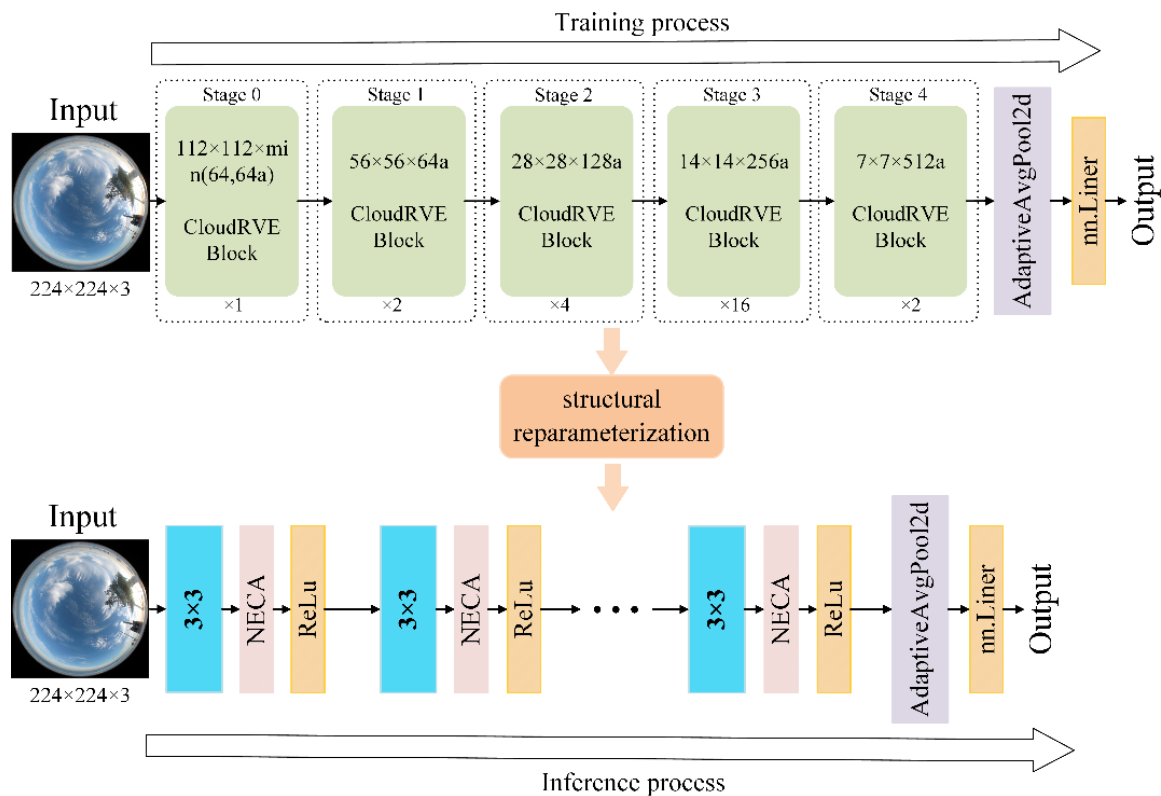
119 CNN plays an important role in the field of target detection, image classification, and image segmentation,
120 especially in the tasks of power line fault detection (Zhao et al., 2016), face recognition (Meng et al., 2021), and
121 medical image segmentation (Zhang et al., 2021), and has been widely used and achieved great achievements. Ground-
122 based cloud image classification is an emerging task in the field of image classification and has achieved rapid and
123 considerable development based on the CNN method. However, it still has some shortcomings such as shallow
124 network level of ground-based cloud image classification method, limited ground-based cloud image classification

125 performance, and small ground-based cloud image classification dataset, which cannot verify the generalization ability
126 of large-scale ground-based cloud image classification dataset.

127 To solve the above problems, the current study improved the RepVGG (Ding et al., 2021) and used it as a basis for
128 elaborating a new classification method for ground-based cloud images called CloudRVE (Cloud Representative
129 Volume Element Network). In this method, the ground-based cloud image was incorporated into the CNN model, and
130 its image features were extracted. Multi-branch convolution layer and channel attention module were used to capture
131 local and global features of the cloud image simultaneously time to enhance the classification performance of ground-
132 based cloud images. The method's application to the multi-modal ground-based cloud dataset named MGCD (Liu et
133 al., 2020a) and ground-based remote sensing cloud database (GRSCD) (Liu et al., 2020b) . The main contributions of
134 this paper are as follows:

- 135 (1) This study elaborated the Improved RepVGG ground-based cloud image classification method with attention
136 convolution called CloudRVE. It broadened the residual structure and comprehensively combined the attention
137 mechanism's abilities to extract the cloud image's global features and describe in detail its local features in the
138 classification process.
- 139 (2) In particular, the Efficient Channel Attention network (ECA) was improved and incorporated into the feature
140 extraction process of ground-based cloud images, which optimization occurred through local cross-channel
141 interaction without dimensionality reduction. Besides, the structural re-parameterization in the inference stage
142 was performed, reducing the model complexity, improving the feature extraction performance, and enhancing
143 the network's learning ability of ground-based cloud image features.
- 144 (3) The comparative analysis of experimental results on the ground-based cloud image classification dataset MGCD
145 proved that the proposed method outperformed ten other state-of-the-art methods in classification accuracy. Its
146 application to GRSCD dataset further verified its generalization ability. Finally, the proposed method's training
147 process optimization and dynamical adjustment of its learning rate were provided by the warm-up method, and
148 the respective recommendations were drawn.

149 The rest of this paper is organized as follows. Section 2 elaborates on the structure and composition of the proposed
150 CloudRVE method for classifying ground cloud images. Section 3 briefly introduces the ground cloud image
151 classification datasets used in this paper and the model evaluation indices. Section 4 provides the experimental results
152 and discusses the feasibility and effectiveness of the proposed method. Finally, Section 5 concludes the study and
153 outlines future research directions and practical application of the research results.



156

157 **Figure 2: CloudRVE network framework. Ground-based cloud images come from Kiel-F datasets (Kalisch and Macke,**
 158 **2008).**

159 This section shows the overall architecture of the proposed RepVGG-based improved classification method, as shown
 160 in Figure 2. In the CloudRVE training process, CloudRVE Block with a multi-branch topology structure is used to
 161 extract features of ground-based cloud images. The multi-branch topology structure has rich gradient information and
 162 a complex network structure, which can effectively improve the characterization ability of local feature information
 163 of ground-based cloud images. Feature maps extracted by CloudRVE Block enter the New Efficient Channel Attention
 164 (NECA) network and learn the feature relationships between sequences to obtain the global feature representation of
 165 an image. In addition, the warm-up method is incorporated into the CloudRVE training process to dynamically
 166 optimize the learning rate and accelerate the model parameter convergence to enhance the model training effect.
 167 CloudRVE inference process uses the single branch topology structure of VGG-style (Simonyan and Zisserman, 2015),
 168 and through structural re-parameterization, the multi-branch convolutional layer and batch normalization (BN) (Ioffe
 169 and Szegedy, 2015) are converted into a 3×3 convolutional layer, increasing its inference speed. The CloudRVE
 170 training process and inference process use the linear classifier to classify the ground-based cloud images to get the

171 final result. The specific framework parameter information of the model is shown in Table 1, where a and b are
 172 magnification factors used to control the network width. The specific contents of each part are as follows.

173 Table 1. The details of CloudRVE training architecture.

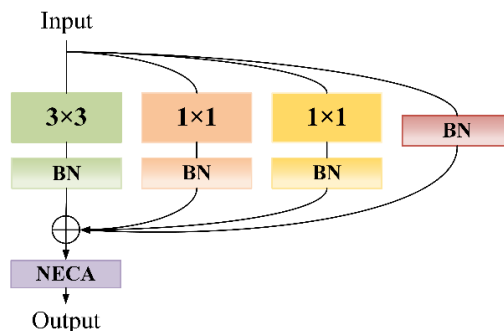
Stage	Blocks of each stage	Output size	Output channels
0	1	224×224	Min (64, 64a)
1	2	112×112	64a
2	4	56×56	128a
3	14	28×28	256a
4	1	14×14	512b

174 2.2 Broadening the CloudRVE Block of Residual Structure

175 CNN is a deep learning model including convolution calculation, including feedforward neural network, which has
 176 representation learning ability, similar to artificial neural network multilayer perceptron (Shi et al., 2017). In 2014,
 177 the most representative convolution neural network VGG came out, which adopted a single-branch topology structure,
 178 greatly improved the image processing effect and model inference speed, and became a new direction for scholars to
 179 learn and develop. With the in-depth study of the VGG, its potential in image processing is close to saturation. Scholars
 180 realize that the VGG has some shortcomings such as simple network structure, few network layers, and large
 181 parameters, which makes it difficult to extract high-order features of images and has limited image-processing
 182 performance. Therefore, improving network complexity and increasing the number of network layers has become a
 183 new research direction. The ResNet developed by (He et al., 2016) differed from the traditional neural network
 184 stacked by convolution layer and pooling layer. The network was stacked by residual modules, which not only
 185 increased the complexity of the network structure and reduced the number of network parameters, but also perfectly
 186 solved the problem of gradient disappearance or gradient explosion caused by increasing the number of network layers,
 187 which could extract abstract image features with semantic information and effectively improve image-processing
 188 performance. By improving the complexity and depth of the network, the ResNet could train the CNN model with
 189 higher accuracy, but there were numerous redundancies in its residual network, impeding the network inference speed
 190 and reducing the accuracy of image processing results (Szegedy et al., 2015). Therefore, increasing the complexity
 191 and depth of the network, weakening its influence on inference speed, and improving the classification effect of
 192 ground-based cloud images become the key goals of this study.

193 To improve the classification effect of the ground-based cloud images, the CloudRVE training process is composed
 194 of CloudRVE blocks that adopt the multi-branch topology. The CloudRVE Block contains four branches and the
 195 improved channel attention module NECA. Its main branch contains a convolutional layer with a convolution kernel
 196 size of 3×3 , which can inspect the input images with a larger neighborhood scope and extract global features easily.

197 Ground-based cloud images contain abundant cloud shape and cloud amount information, while a large convolution
 198 kernel tends to ignore cloud boundary features, resulting in inadequate feature extraction from ground-based cloud
 199 images. Therefore, the two bypass branches of CloudRVE Block adopt the convolution layer with the convolution
 200 kernel size of 1×1 , which can not only extract fine cloud boundary features and abstract cloud cover features but also
 201 keep the output dimension consistent with the input dimension, facilitating the multi-branch ground-based cloud
 202 image feature fusion. The third bypass branch of CloudRVE Block adopts the Identity branch, whose purpose is to
 203 take the input as the output and change the learning objective to the residual result approaching 0 so that the accuracy
 204 does not decline with the deepening of the network. In addition, each branch is connected to the BN layer, not only to
 205 avoid overfitting but also to prevent gradient disappearance or explosion. The specific structure of CloudRVE Block
 206 is shown in Figure 3. The input feature maps pass through three branches with a convolutional layer and BN layer at
 207 the same time. The output obtained by the input feature maps is summed with the Identity branch and input into the
 208 NECA module to obtain the final output feature.



209

210 **Figure 3: CloudRVE Block structure.**

211 2.3 NECA Module Focusing on Full Image Features

212 The attention mechanism is to let the neural network have the information processing way to distinguish the key points
 213 and to capture the connection between global information and local information flexibly. Its purpose is to enable the
 214 model to obtain the target region that needs to be focused on, put more weight on this part, highlight significant useful
 215 features, and suppress and ignore irrelevant features. The NECA (New Efficient Channel Attention) is an
 216 implementation form of channel attention mechanism, which can strengthen channel features without changing the
 217 size of the input feature maps. It adopts a local cross-channel interaction strategy without dimensionality reduction so
 218 that the 1×1 convolution layer can replace the full connection layer to learn channel attention information, which can
 219 effectively avoid the negative impact of dimensionality reduction on channel attention learning. The network
 220 performance is guaranteed and the complexity of the model is significantly reduced.

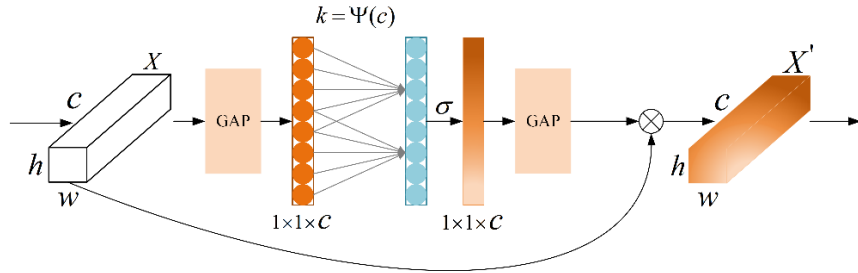
221 The ground-based cloud image samples in Figure 2 were taken by the all-sky imager and could cover the sky in
 222 this area. However, the ground-based cloud images contain not only the valid area of the whole sky but also the black

223 invalid area. Therefore, the NECA module abandons the traditional global maximum pooling and adopts double global
 224 average pooling. The global average pooling formulas are as follows:

$$225 \quad \gamma_{gap} = \frac{1}{wh} \sum_{i=1, j=1}^{w, h} X_{ij}, \quad X \in R^{w \times h \times c}, \quad (1)$$

$$226 \quad \eta_{gap} = \sigma(V_k^{gap} \gamma_{gap}), \quad V_k^{gap} \in R^{c \times c}, \quad (2)$$

227 where X and X' represent the input and output feature maps, respectively, whereas w , h , and c are the width, height,
 228 and number of channels of the input feature map. The NECA module adopts a double global average pool, which can
 229 effectively improve its noise suppression ability and enhance its channel feature extraction ability, which can avoid
 230 the black invalid part of the feature calculation. The NECA module structure is shown in Figure 4.



231

N

232 **Figure 4: NECA model structure.**

233 Here b and r are fixed values, and their values are set to 1 and 2, respectively, while k represents the convolution
 234 kernel size and has a corresponding relationship with c . As the network deepens, the number of channels c increases
 235 by the power of 2. Therefore, k should not be a fixed value, but a dynamic change and its relationship are as follows:

$$236 \quad C = \phi(k) = 2^{(\gamma * k - b)} \quad (3)$$

$$237 \quad K = \psi(C) = \left\lfloor \frac{\log_2(c)}{r} - \frac{b}{r} \right\rfloor_{odd} \quad (4)$$

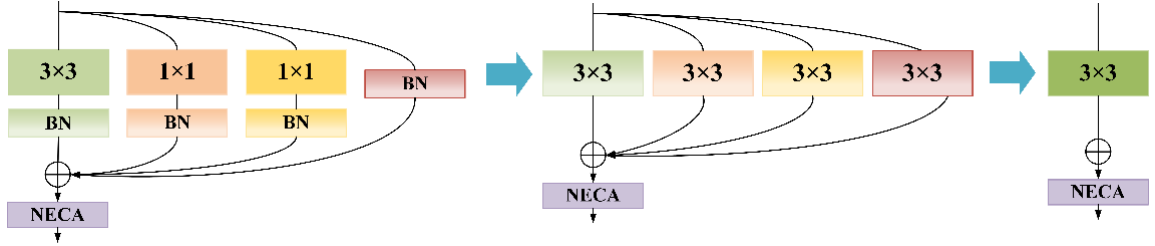
238 2.4 Inference Process from Multi-Branch to Single-Branch

239 The residual module is crucial to the CloudRVE training process. Its multi-branch topology can improve CloudRVE
 240 Block's ability to extract ground cloud image features and solve optimization problems such as gradient disappearance
 241 and gradient explosion caused by increasing network depth. However, the multi-branch topology will occupy more
 242 memory for the CloudRVE reasoning process, resulting in insufficient utilization of hardware computing power and
 243 slower reasoning speed. If the single-branch topology is adopted, the computing load is reduced and the inference
 244 time is saved, thus reducing memory consumption. Therefore, the single-branch topology structure is adopted in the
 245 CloudRVE inference stage, and the trained CloudRVE Block needs to be transformed into a single-branch topology
 246 model through structural re-parameterization. The conversion process mainly includes the fusion of the convolutional

247 layer and BN layer, the conversion of the BN layer into a convolutional layer, and the fusion of the multi-branch
 248 convolutional layer. We use $W_{(3)} \in R^{C_1 \times C_2 \times 3 \times 3}$ as 3×3 convolution layers, and use C_1, C_2 as input channels and
 249 output channels respectively, and use $W_{(1)} \in R^{C_1 \times C_2 \times 1 \times 1}$ as 1×1 convolution layers. In addition, we use $\mu_{(3)}, \sigma_{(3)},$
 250 $\gamma_{(3)}, \beta_{(3)}$ to represent the mean value, standard deviation, learning scaling factor, and deviation of the BN layer of the
 251 main branch, and use $\mu_{(1)}, \sigma_{(1)}, \gamma_{(1)}, \beta_{(1)}$ to represent the parameters of the BN layer of the by-pass branch containing
 252 1×1 convolution layer, and use $\mu_{(0)}, \sigma_{(0)}, \gamma_{(0)}, \beta_{(0)}$ to represent the parameters of the BN layer of the identity branch,
 253 and use $M_{(1)} \in R^{N \times C_1 \times H_1 \times W_1}, M_{(2)} \in R^{N \times C_2 \times H_2 \times W_2}$ to represent the input and output. The CloudRVE Block structure
 254 reparameterization calculation process is as follows:

$$255 \quad M_{(2)} = \text{BN}(M_{(1)} * W_{(3)}, \mu_{(3)}, \sigma_{(3)}, \gamma_{(3)}, \beta_{(3)}) + \text{BN}(M_{(1)} * W_{(1)}, \mu_{(1)}, \sigma_{(1)}, \gamma_{(1)}, \beta_{(1)}) \\ + \text{BN}(M_{(1)} * W_{(1)}, \mu_{(1)}, \sigma_{(1)}, \gamma_{(1)}, \beta_{(1)}) + \text{BN}(M_{(1)}, \mu_{(0)}, \sigma_{(0)}, \gamma_{(0)}, \beta_{(0)}) \quad (5)$$

256 The input feature map is inputted into the NECA module through the 3×3 convolution layer completed by fusion.
 257 The process is shown in Figure 5.



258
 259 **Figure 5: Re-parameterization process of CloudRVE Block structure.**

260 2.4.1 Fusion of Convolutional Layer and BN Layer

261 This section first describes the fusion of the main branch 3×3 convolution layer with the BN layer and then describes
 262 the transformation of the bypass branch 1×1 convolution layer into the 3×3 convolution layer and fusion with the BN
 263 layer. In the inference stage, the number of convolutional kernel channels in the convolution layer is the same as the
 264 number of channels in the input feature map, and the number of convolutional kernel channels in the output feature
 265 map is the same. The main parameters of the BN layer include mean μ , variance σ^2 , learning ratio factor γ , and
 266 deviation β . Of these, μ and σ^2 are obtained statistically in the training stage, while γ and β are obtained by learning
 267 in the training stage. The calculation of the i channel of the input BN layer is performed as follows:

$$268 \quad y_i = \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}} \times \gamma_i + \beta_i, \quad (6)$$

269 where x is the input and ϵ is the constant approaching 0. The calculation process of the i channel input BN in the
 270 feature map can be expressed as follows:

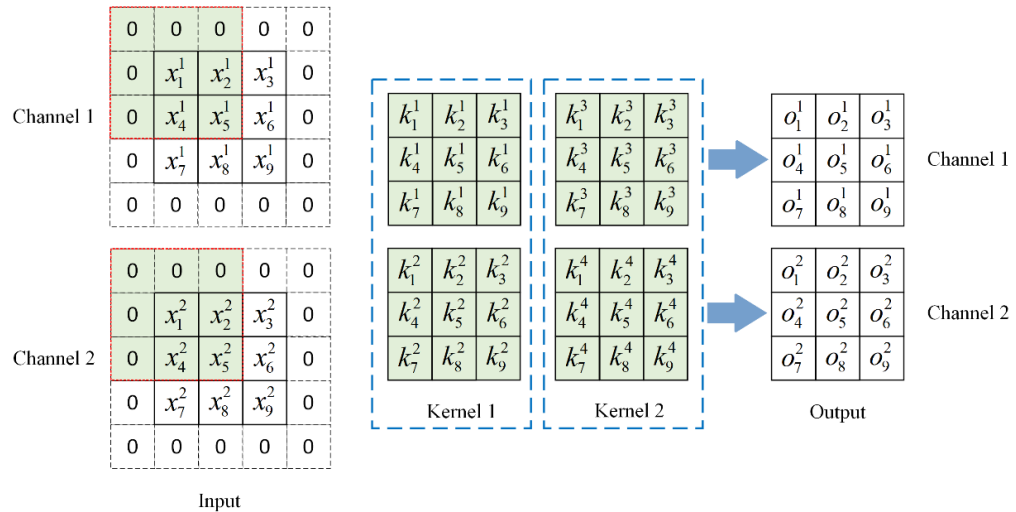
$$271 \quad bn(M, \mu, \sigma, \gamma, \beta)_{:,i,:} = (M_{:,i,:} - \mu_i) \frac{\gamma_i}{\sigma_i} + \beta_i = \frac{\gamma_i}{\sigma_i} M_{:,i,:} + \beta_i - \frac{\gamma_i}{\sigma_i} \mu_i, \quad (7)$$

272 where M is the output feature map obtained by weighted summation of the convolution layer, input to BN layer and
 273 ignore x . Therefore, we can multiply γ_i/σ_i to the i convolution kernel of the 3×3 convolution layer:

$$274 \quad W'_{i,:} = \frac{\gamma_i}{\sigma_i} W_{i,:} \quad (8)$$

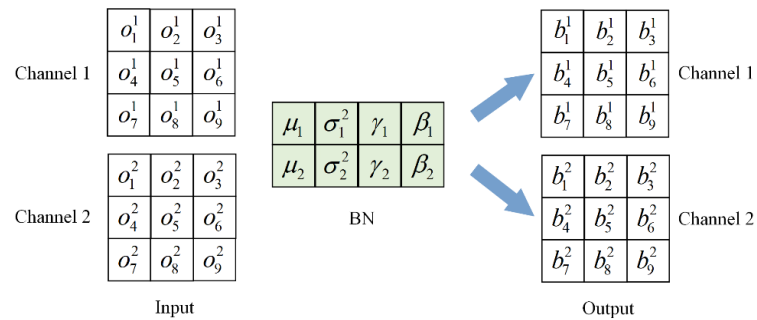
$$275 \quad b'_i = \beta_i - \frac{\mu_i \gamma_i}{\sigma_i} \quad (9)$$

276 The i convolution kernel weight of the fusion of the 3×3 convolution layer and BN layer is obtained, and the
 277 specific fusion process is shown in Figures 6 and 7. The input channel C_1 and output channel C_2 make two, and the
 278 stride is one. In the convolution layer, the input feature map is calculated by convolution to obtain the output feature
 279 map with the number of channels 2. Figure 8 shows that the number of channels in the BN layer is 2, and the output
 280 feature map of the convolution layer is used as the input feature map of the BN layer. The output feature map with
 281 the number of channels being 2 is obtained via equation (2).



282

283 **Figure 6: Input feature map through convolution layer process. For visualization, we assume that $C_1=C_2=2$.**



284

285 **Figure 7: Convolutional layer output feature map through the BN layer process.**

286 In addition, to ensure that the size of the output feature map is consistent with that of the input feature map, the input
 287 feature map should be converted into 5×5 size by padding operation. The concrete convolution is as follows:

$$288 \quad o_1^1 = x_1^1 \cdot k_5^1 + x_2^1 \cdot k_6^1 + x_4^1 \cdot k_8^1 + x_5^1 \cdot k_9^1 + x_1^2 \cdot k_5^2 + x_2^2 \cdot k_6^2 + x_4^2 \cdot k_8^2 + x_5^2 \cdot k_9^2 \quad (10)$$

289 The specific calculation process of the input feature map through the BN layer is

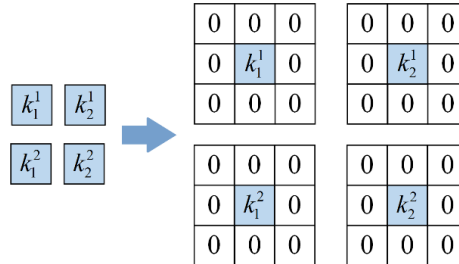
$$290 \quad b_1 = \frac{(x_1^1 \cdot k_5^1 + x_2^1 \cdot k_6^1 + x_4^1 \cdot k_8^1 + x_5^1 \cdot k_9^1 + x_1^2 \cdot k_5^2 + x_2^2 \cdot k_6^2 + x_4^2 \cdot k_8^2 + x_5^2 \cdot k_9^2) - \mu_1}{\sqrt{\sigma^2 + \varepsilon}} \cdot \gamma_1 + \beta_1 \quad (11)$$

291 Re-arranging equation (7) yields

$$292 \quad b_1 = (x_1^1 \cdot k_5^1 + x_2^1 \cdot k_6^1 + x_4^1 \cdot k_8^1 + x_5^1 \cdot k_9^1 + x_1^2 \cdot k_5^2 + x_2^2 \cdot k_6^2 + x_4^2 \cdot k_8^2 + x_5^2 \cdot k_9^2) \cdot \frac{\gamma_1}{\sqrt{\sigma^2 + \varepsilon}} + (\beta_1 - \frac{\mu_1}{\sqrt{\sigma^2 + \varepsilon}}) \quad (12)$$

$$293 \quad c = \frac{\gamma_1}{\sqrt{\sigma^2 + \varepsilon}} ; \quad d = \beta_1 - \frac{\mu_1}{\sqrt{\sigma^2 + \varepsilon}} \quad (13)$$

294 In equation (8), c and d are constants and are multiplied to the first convolution kernel of the convolution layer to
 295 obtain the parameters of the first convolution kernel after the convolution layer and BN layer are fused. Other fused
 296 convolution kernel parameters are calculated similarly. The convolution layer and BN layer are fused by the bypass
 297 branch containing a 1×1 convolution layer. The convolution layer is first converted to 3×3 size by padding operation
 298 and then fused with the BN layer by repeating the above steps. The convolution layer padding process is shown in
 299 Figure 8.



300
 301 **Figure 8: 1×1 convolution layer transformed into 3×3 convolution layer.**

302 2.4.2 Converting the BN Layer to the Convolution Layer

303 The identity bypass branch has only a BN layer, its function is to ensure the identity mapping of the input feature map
 304 and output feature map. To realize the identical mapping between the input feature map and the output feature map in
 305 the fusion process, a 3×3 convolution layer with 2 convolution kernels and 2 convolution kernel channels needs to
 306 be designed. Secondly, the input feature map needs to be converted into a 5×5 feature map by padding operation.
 307 The specific process is shown in Figure 9. The output feature map is obtained by convolution calculation of the input
 308 feature map, and its parameters and sizes are consistent with those of the input feature map. Finally, the fusion process
 309 of the 3×3 convolution layer and BN layer is repeated to obtain a new 3×3 convolution layer.

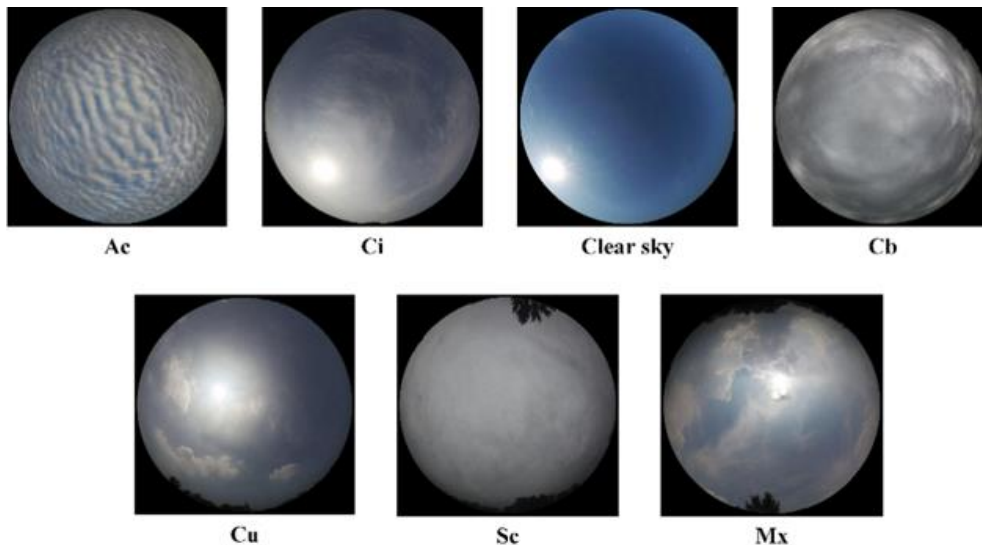
328 3. Dataset and Experimental Settings

329 This section introduces two kinds of ground-based cloud image classification datasets, MGCD and GRSCD, and
330 describes the relevant experimental Settings. Subsection 3.1 describes MGCD and GRSCD in detail, and Subsection
331 3.2 details experimental setting parameters and model evaluation indices.

332 3.1 Ground-Based Cloud Image Dataset

333 3.1.1 Introduction to MGCD Dataset

334 Multi-modal Ground-based Cloud image Dataset (MGCD) is the first ground-based cloud image classification dataset
335 composed of ground-based cloud images and multi-modal information, which was collected by the School of
336 Electronics and Communication Engineering of Tianjin Normal University and the Meteorological Observation
337 Center of Beijing Meteorological Bureau of China from 2017 to 2018. There are 8000 ground-based cloud images in
338 MGCD, and 4000 ground-based cloud images in the training set and testing set, including altocumulus (Ac), cirrus
339 (Ci), clear sky (Cl), cumulonimbus (Cb), cumulus (Cu), stratocumulus (Sc), and mix (Mx). In addition, cloud images
340 with a cloud cover of less than 10% are classified as clear sky, and each sample contains a captured ground cloud
341 image and a set of multimodal cloud information. Among them, the ground-based cloud images are collected by an
342 all-sky camera with a fisheye lens, and its data storage format is JPEG with a resolution of 1024×1024 pixels;
343 Multimodal information is collected by weather stations, including temperature, humidity, pressure, and wind speed,
344 and these four elements are stored in the same vector. Figure 10 is a partial sample of the MGCD dataset, and the
345 specific information is shown in Table 2.



346

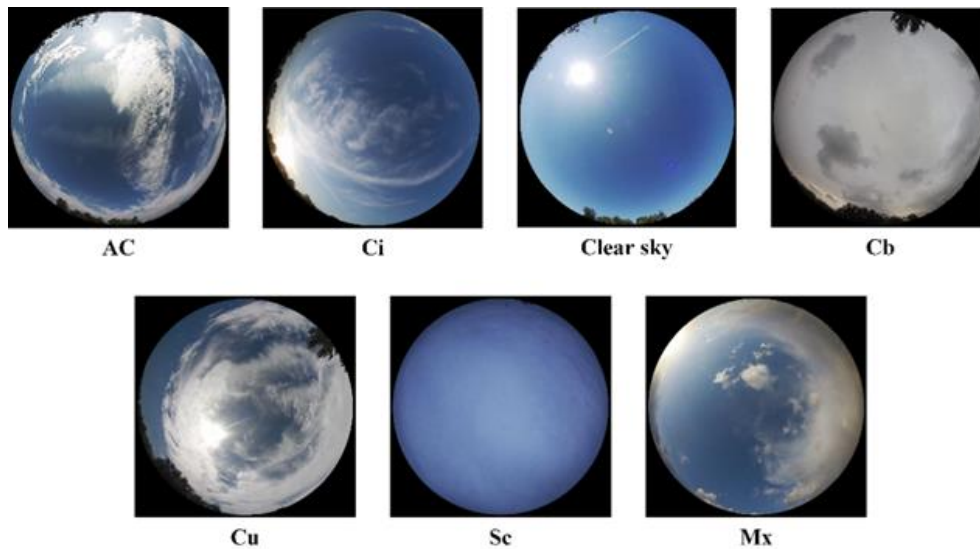
347 **Figure 10: Sample legend of MGCD dataset (Liu et al., 2020a).**

Table 2. MGCD dataset-specific information.

No	Class	Training	Testing	Total
1	Ac	365	366	731
2	Ci	662	661	1323
3	Cl	669	669	1338
4	Cb	593	594	1187
5	Cu	719	719	1438
6	Sc	482	481	963
7	Mx	510	510	1020
Total		4000	4000	8000

349 3.1.2 Introduction to GRSCD Dataset

350 Ground remote sensing cloud dataset (GRSCD) is a ground-based cloud image classification dataset composed of
351 ground-based cloud images and multimodal information. It was collected by the College of Electronic and
352 Communication Engineering of Tianjin Normal University and the Meteorological Observation Center of Beijing
353 Meteorological Administration of China from 2017 to 2018. The total number of ground-based cloud images in
354 GRSCD is consistent with MGCD, with a training set and a testing set each accounting for 50%, including 7 types of
355 clouds: altostratus (Ac), cirrus(Ci), clear sky(Cl), cumulonimbus(Cb), cumulus(Cu), stratocumulus(Sc), and mix(Mx).
356 Among them, the features of cumulonimbus and stratocumulus in MGCD are not distinct and easy to confuse; the
357 features of altostratus and cumulus in GRSCD are not distinct and easy to confuse. In addition, each sample contains
358 a ground-based cloud image and a set of multi-modal cloud information, and cloud images with cloud cover not
359 exceeding 10% are classified as clear sky. Figure 11 depicts a partial sample of the GRSCD dataset. The specific data
360 are listed in Table 3.



361

362 **Figure 11: Sample legend of GRSCD dataset (Liu et al., 2020b).**

Table 3. GRSCD dataset-specific information

No	Class	Training	Testing	Total
1	Ac	400	331	731
2	Ci	650	673	1323
3	Cl	650	688	1338
4	Cb	600	587	1187
5	Cu	690	748	1438
6	Sc	500	463	963
7	Mx	510	510	1020
	Total	4000	4000	8000

364 3.2 Experimental Setting

365 3.2.1 Implementation Details

366 All experiments in this paper adopt Python programming language and run on Intel(R) Core (TM) i9-12700K CPU
367 @ 3.60GHz. NVIDIA GeForce RTX 3090 24G Graphical Processing Unit (GPU) platform and uses Pytorch as a deep
368 learning framework. The CNN experiment is trained on the ground-based cloud image classification datasets MGCD
369 and GRSCD respectively. The number of training data accounts for 50%, the initial learning rate is set to 0.0002,
370 Batchsize is set to 32, and Adam optimizer (Kingma and Ba, 2015) is used to optimize all available parameters in the
371 network. In addition, to improve the generalization ability of the CNN model and the convergence speed of the
372 experiment, the transfer learning method is adopted in the training stage, and model parameters are obtained by
373 training RepVGG with the ground-based cloud image classification dataset made by the team and used as the weight
374 of pre-training. CNN experiment directly trains based on pre-training weight, which can accelerate the model
375 convergence speed and shorten the training time, avoid the problem of parameter overfitting, and promote the rapid
376 gradient decline.

377 3.2.2 Evaluation Index

378 To objectively evaluate the ground-based cloud image classification performance of CloudRVE and other CNN
379 models, the accuracy rate, recall rate, and the average values of different indices of 7 types of clouds in MGCD and
380 GRSCD datasets are calculated in the experiment, which is used as evaluation indices of CNN model. The accuracy
381 rate and average accuracy rate is derived based on positive and negative samples, n represents the number of cloud
382 types and the calculation process is as follows:

$$383 \text{ Accuracy}(Acc) = \frac{TP+TN}{TP+TN+FP+FN}, \overline{\text{Accuracy}}(\overline{Acc}) = \frac{1}{n} \sum_{i=1}^n \frac{TP_i+TN_i}{TP_i+TN_i+FP_i+FN_i} \quad (15)$$

384 TP (True Positive) parameter is the number of correctly classified samples for a specific genus, TN (True Negative)
385 parameter is the number of correctly classified samples for the remaining genus, and FN (False Negative) parameter
386 is the number of misclassified samples for a specific class genus. FP (False Positive) parameter is the number of

387 misclassified samples for the remaining classes genera. The precision rate, average precision rate, recall rate and
 388 average recall rate can be expressed as:

$$389 \text{ Precision}(Pr) = \frac{TP}{TP+FP}, \overline{\text{Precision}}(\overline{Pr}) = \frac{1}{n} \sum_{i=1}^n \frac{TP_i}{TP_i+FP_i} \quad (16)$$

$$390 \text{ Recall}(Re) = \frac{TP}{TP+FN}, \overline{\text{Recall}}(\overline{Re}) = \frac{1}{n} \sum_{i=1}^n \frac{TP_i}{TP_i+FN_i} \quad (17)$$

391 In addition, the specificity, average specificity, F1_score and average F1_score are also used as evaluation indices
 392 of the CNN model in the experiment, and their expressions are shown as follows:

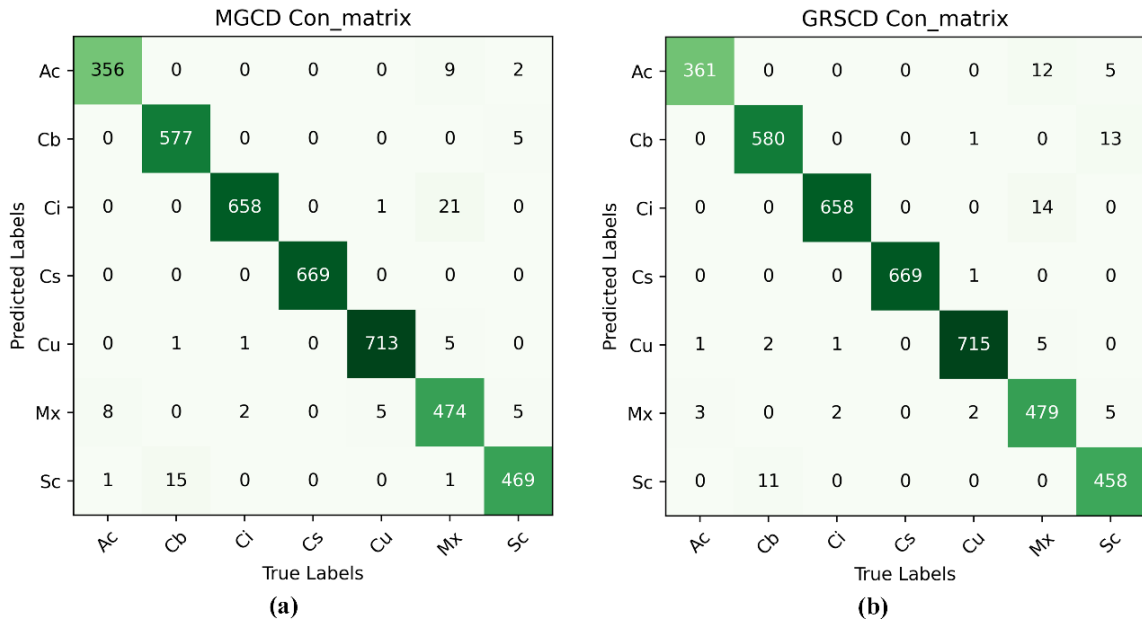
$$393 \text{ Specificity}(TNR) = \frac{TN}{FP+TN}, \overline{\text{Specificity}}(\overline{TNR}) = \frac{1}{n} \sum_{i=1}^n \frac{TN_i}{FP_i+TN_i} \quad (18)$$

$$394 \text{ F1_score}(F1) = \frac{2 \times Pr \times Re}{Pr+Re}, \overline{\text{F1_score}}(\overline{F1}) = \frac{1}{n} \sum_{i=1}^n \frac{2 \times Pr_i \times Re_i}{Pr_i+Re_i} \quad (19)$$

395 4. Experimental Results and Discussion

396 4.1 Classification Results of Ground-Based Cloud Images

397 Figure 12 shows the confusion matrix of MGCD and GRSCD datasets, showing CloudRVE prediction results on
 398 MGCD and GRSCD datasets. The horizontal axis represents the true cloud image classification, while the vertical
 399 axis represents the predicted cloud image classification, where the value of the diagonal element represents the correct
 400 number of cloud image classifications and the value of the off-diagonal element represents the number of cloud image
 401 classification errors. As can be seen from Figure 12(a), in the MGCD dataset, the correct classification of the Cu is
 402 the largest, while the misclassification of the cloud images mainly comes from Sc and Mx. The reason is that the cloud
 403 base of Sc is blackened by illumination, making it easily confused with Cb. In addition, the dynamic change of cloud
 404 will lead to a change in the viewpoint of the whole sky camera, thus increasing the difficulty of cloud genus
 405 identification. As can be seen in Figure 12(b), in the GRSCD dataset, the correctly classified cloud images of the same
 406 Cu had the largest number, while the incorrectly classified ones mainly came from Mx and Sc. The Mx cloud is a
 407 hybrid cloud, containing a variety of different cloud genera, with large shares of Ac, Ci, and Cu, which could be
 408 erroneously classified as Mx. Similarly, Sc could be taken for Cb, due to their similar features, impeding the correct
 409 identification.



410

411 **Figure 12: Confusion matrix images. (a)MGCD confusion matrix image. (b) GRSCD confusion matrix image.**

412 The overall classification accuracy of the CloudRVE method proposed in this paper in MGCD and GRSCD datasets
 413 and the classification results of each cloud genus are listed in Tables 4 and 5. It can be seen that the accuracy of
 414 CloudRVE in MGCD and GRSCD datasets reached 98.15 and 98.07%, respectively. The characteristics of the Ci in
 415 MGCD and GRSCD datasets were easy to identify, resulting in the accuracy rate, recall rate, specificity, and F1 value
 416 reaching 100%. In the MGCD dataset, the accuracy rate, recall rate, and F1 value of the other six cloud genera all
 417 exceeded 95.00%, and the specificity was above 99.50%. The accuracy and specificity of the Ci were the highest,
 418 reaching 98.64 and 99.73%, respectively. Cu had the highest recall rate and F1 value, reaching 99.17 and 98.89%,
 419 respectively. In addition, the recall rate and F1 value of Sc and Mx were about 2.00% lower than other cloud genera,
 420 mainly their characteristics in the MGCD dataset were similar to those of Cb and Ci, respectively, reducing
 421 CloudRVE's ability to classify them.

422

Table 4. Classification results for the MGCD dataset.

Genus	\overline{Acc} (%)	Pr (%)	Re (%)	TNR (%)	F1 (%)
Cu		98.62	99.17	99.70	98.89
Ac		97.02	98.08	99.70	97.55
Ci		98.64	98.94	99.73	98.79
Cl	98.15	100.0	100.0	100.0	100.0
Sc		97.26	95.84	99.63	96.54
Cb		97.13	97.13	99.51	97.13
Mx		97.24	96.67	99.60	96.95

Table 5. Classification results for the GRSCD dataset.

Genus	\overline{Acc} (%)	Pr (%)	Re (%)	TNR (%)	F1 (%)
Cu		99.30	99.03	99.85	99.16
Ac		94.24	98.63	99.39	96.39
Ci		97.91	99.24	99.58	98.57
Cl	98.07	100.0	100.0	100.0	100.0
Sc		98.10	96.47	99.74	97.27
Cb		97.33	98.48	99.53	97.90
Mx		97.74	93.33	99.68	95.49

424 In the GRSCD dataset, the accuracy rate, recall rate, and F1 value of the other six cloud genera exceeded 94.00%,
 425 and the specificity as over 99.30%. Cu had the highest accuracy, specificity, and F1 value, reaching 99.30, 99.85, and
 426 99.16%. The recall rate of Ci was the highest, reaching 99.17%. In addition, the Ac accuracy was only 94.24%, mainly
 427 because Ac contained a small amount of Sc, and CloudRVE could easily to misjudge Ac as Sc or Mx. Mx contained
 428 a variety of other clouds, and the images composition was complex. Cloud clusters of different can genera varied in
 429 size and shape, resulting in lower recall rate and F1 values.

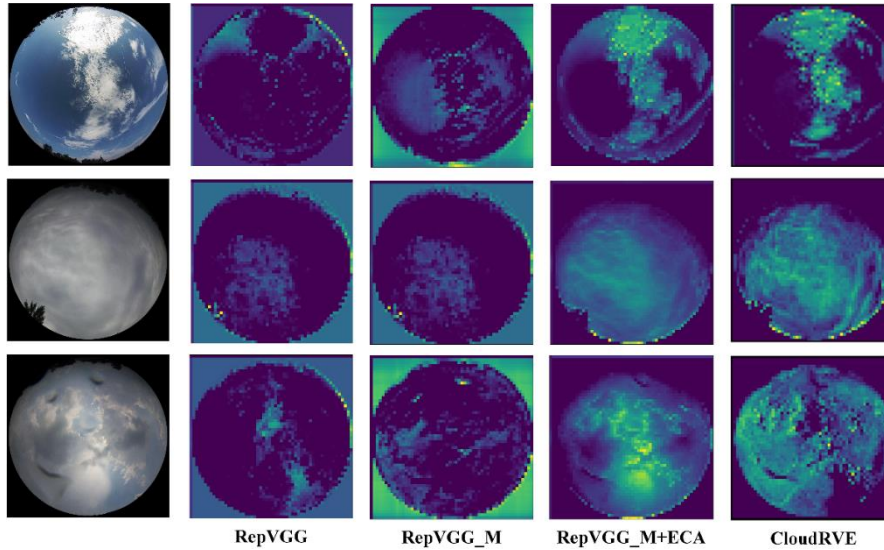
430 4.2. Ablation Experiment

431 **Table 6.** Results of the ablation experiment.

Dataset	Model	\overline{Acc} (%)	\overline{Pr} (%)	\overline{Re} (%)	\overline{TNR} (%)	$\overline{F1}$ (%)
MGCD	RepVGG	95.57	95.31	94.99	99.26	95.14
	RepVGG_M	95.97	95.65	95.67	99.33	95.56
	RepVGG_M+ECA	96.80	96.60	96.37	99.47	96.45
	CloudRVE	98.15	97.99	97.98	99.68	97.83
GRSCD	RepVGG	95.42	94.99	94.88	99.24	94.92
	RepVGG_M	95.70	95.46	95.30	99.29	95.36
	RepVGG_M+ECA	96.10	95.67	95.74	99.35	95.68
	CloudRVE	98.07	97.80	97.88	99.68	97.82

432 In this section, the ablation experiment is used to compare the original structure and different improvement stages of
 433 the proposed method on the MGCD and GRSCD datasets respectively, and the results are shown in Table 6.
 434 RepVGG_M is the main improved network, ECA is the attention module, CloudRVE is the combined improved
 435 network of RepVGG_M and NECA, and is the final version of the method proposed in this paper. It can be seen from
 436 the data in the table that the performance of each improvement stage of the network model is improved compared to
 437 the previous stage, which not only verifies the feasibility of extracting more cloud image detail features by adding
 438 1×1 convolutional layer branches but also verifies that NECA can effectively improve the noise suppression ability
 439 and enhance the channel feature extraction ability. Compared with the original network structure, the accuracy of
 440 CloudRVE in the MGCD dataset increased by 2.58%, the average accuracy rate increased by 2.68%, the average
 441 recall rate increased by 2.99%, the average specificity increased by 0.42%, and the average F1 value increased by

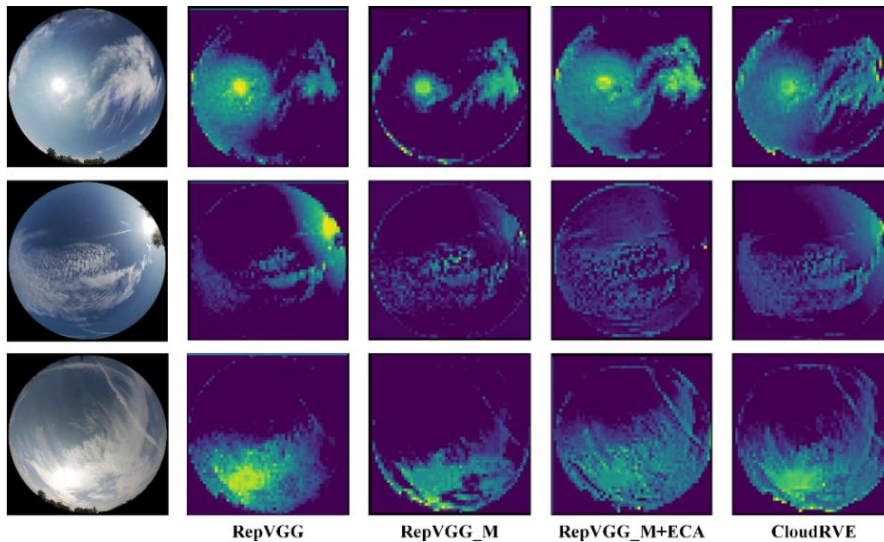
442 2.69%. In the GRSCD dataset, the accuracy rate increased by 2.65%, the average accuracy rate increased by 2.81%,
 443 the average specificity increased by 0.44%, and the average F1 value increased by 2.69%. Therefore, it can be seen
 444 from the data display that the method proposed in this paper has the best performance.



445

446 **Figure 13: Feature extraction of different models based on MGCD (Liu et al., 2020a).**

447



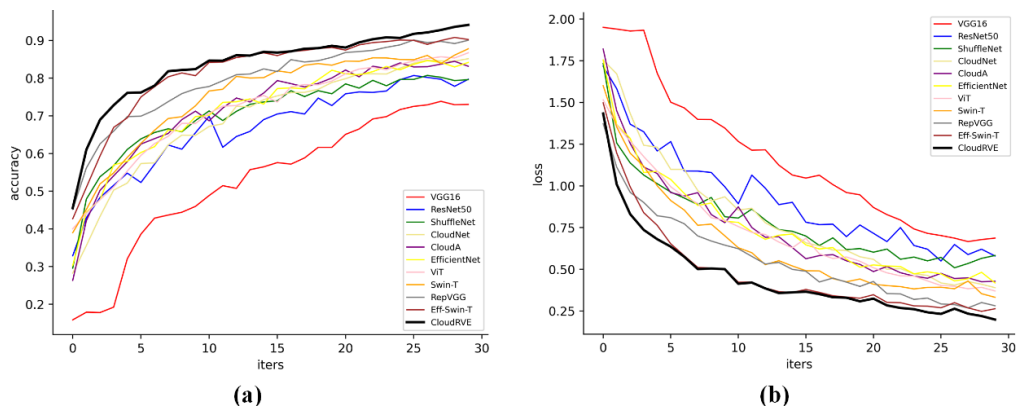
448

449 **Figure 14: Feature extraction of different models based on GRSCD (Liu et al., 2020b).**

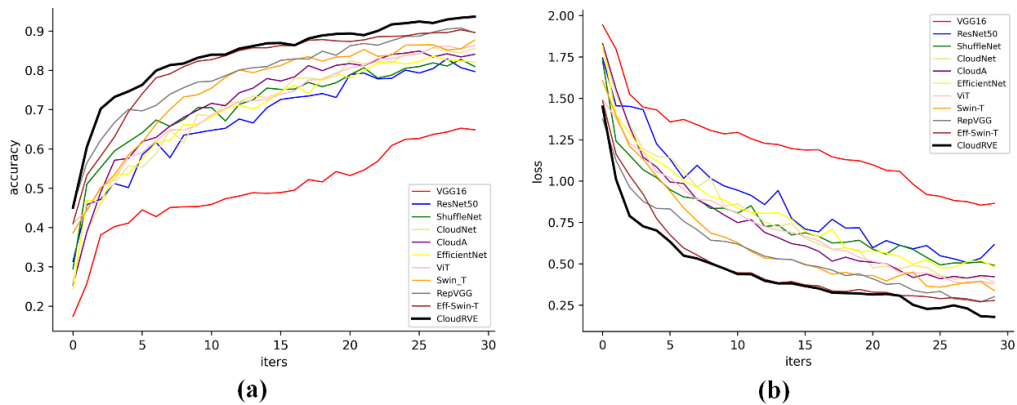
450 To visually compare the performance of the original structure and the method proposed in this paper in different
 451 improvement stages, we visualize the features by extracting the feature map of the middle layer of the network and

452 then explain the feature extraction ability of the original structure and the method proposed in this paper in different
 453 improvement stages, as shown in Figures 13 and 14. The method generates a rough feature map to display the
 454 important region of the predicted images through the parameter weights generated by network training, in which the
 455 brighter the region indicates the higher its importance, and the darker the region represents the sky or those that cannot
 456 be extracted. Figure 13 shows that CloudRVE has the best feature location and extraction ability by showing the
 457 feature maps of three different cloud images in the MGCD dataset. Figure 14 shows that the three cloud images of the
 458 GRSCD dataset include not only clouds and sky but also strong sunlight, which affects the classification accuracy of
 459 the model. However, it can be seen from the feature maps that CloudRVE not only has the best feature extraction
 460 ability but also has a strong ability to suppress noise such as sunlight.

461 4.3 Comparison of Experimental Results



462 (a) (b)
 463 **Figure 15: Training accuracy (a) and training loss (b) curves of the MGCD dataset.**



464 (a) (b)
 465 (a) (b)
 466 **Figure 16: Training accuracy (a) and training loss (b) curves of the GRSCD dataset.**

467 To verify the feasibility of the proposed CloudRVE method, we compared it with other advanced methods, including
468 CloudNet (Zhang et al., 2018), CloudA (Wang et al., 2020), Eff-Swin-T (Li et al., 2022), and other ground-based
469 cloud image classification methods. These included such classic CNN models as VGG16 (Szegedy et al., 2015),
470 ResNet50 (He et al., 2016), ShuffleNet (Zhang et al., 2018) and EfficientNet (Tan and Le, 2019). In addition, we
471 compared it with other Transformer-based classification models such as ViT-L (Dosovitskiy et al., 2022), Swin-T(Liu
472 et al., 2021), etc. Figures 15 and 16 illustrate the performances of different methods by displaying the training accuracy
473 and training loss curves of MGCD and GRSCD datasets. Here the black bold curve represents the CloudRVE method,
474 which has the largest accuracy value, the fastest convergence rate, the smallest loss rate, and the fastest decline rate
475 in the training stage. This strongly indicates that the CloudRVE method has the best classification performance of
476 ground-based cloud images.

477 **Table 7.** Comparison of experimental results.

Method	MGCD					GRSCD				
	\overline{Acc} (%)	\overline{Pr} (%)	\overline{Re} (%)	\overline{TNR} (%)	$\overline{F1}$ (%)	\overline{Acc} (%)	\overline{Pr} (%)	\overline{Re} (%)	\overline{TNR} (%)	$\overline{F1}$ (%)
VGG-16	78.25	77.04	75.52	96.36	75.55	73.50	73.88	70.29	95.53	70.87
ResNet-50	85.98	85.24	84.55	97.67	84.82	86.51	85.56	85.38	97.75	85.34
ShuffleNet	86.95	86.08	85.68	97.83	85.71	86.99	86.85	85.18	97.82	85.71
CloudNet	90.01	89.24	89.08	98.34	89.13	89.60	89.06	88.60	98.27	88.79
CloudA	89.62	88.78	88.50	98.28	88.61	90.03	89.54	88.71	98.34	89.03
EfficientNet	91.17	90.66	90.22	98.53	90.27	90.10	89.68	88.92	98.35	89.13
ViT-L	91.11	90.91	90.21	98.55	90.40	90.98	90.49	90.33	98.50	90.39
Swin-T	92.87	92.44	91.63	98.63	91.76	93.55	93.22	92.87	98.93	92.71
RepVGG	95.57	95.31	94.99	99.26	95.14	95.42	94.99	94.88	99.24	94.92
Eff-Swin-T	96.93	96.73	96.44	99.49	96.56	95.62	95.41	95.11	99.27	95.21
CloudRVE	98.15	97.99	97.98	99.68	97.83	98.07	97.80	97.88	99.68	97.82

478 The comparative analysis results of the above methods are summarized in Table 7. It can be seen from the
479 experimental results that RepVGG had the best performance among the CNN-based methods. Among them, the
480 accuracy rate has the most significant improvement, and the precision and recall rates also have good improvement.
481 The accuracy rate, precision rate, recall rate for the MGCD dataset reached 95.57, 95.31, and 94.99, respectively, while
482 those for the GRSCD dataset were 95.42, 94.99, and 94.88, respectively. Ground-based cloud images have more
483 texture features and deep semantic features than other images, and more image features need to be obtained to meet
484 the classification requirements of such images. In recent years, Transformer has been widely used for image
485 processing tasks due to its strong feature extraction capability. Several scholars have improved the Transformer
486 derivative model through continuous exploration. Among them, Eff-Swin-T was an improvement based on Swin-T,
487 and its performance on MGCD and GRSCD datasets was better than that of the classic CNN model. Its accuracy rate,
488 precision rate, and recall rate reached 96.93, 96.73, 96.44, and 95.62, 95.41, 95.11, respectively. Compared with

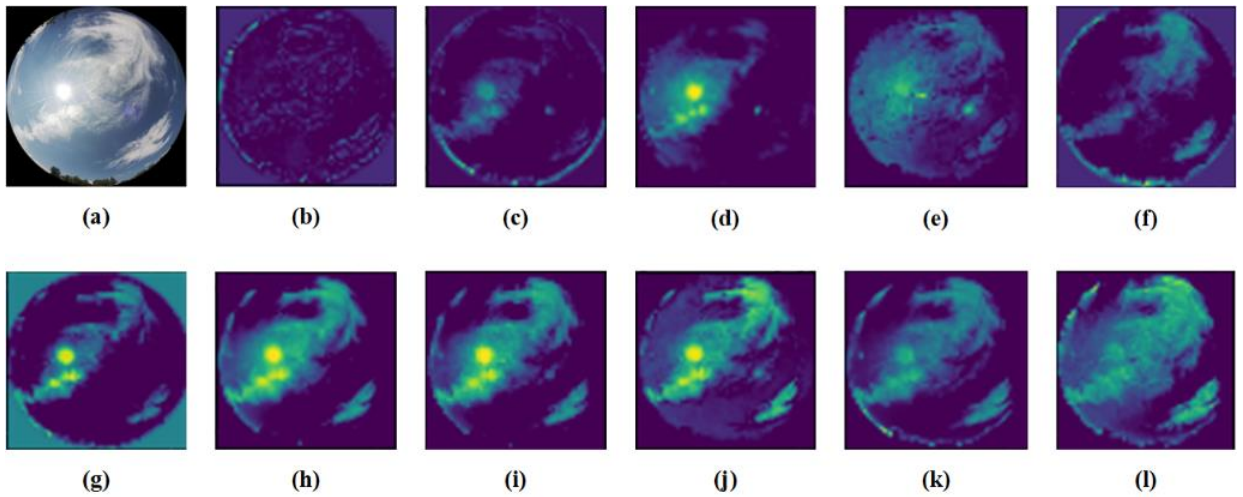
489 Transformer and classical networks, the proposed method had much better classification performance of ground-based
 490 cloud images. For different cloud image classification datasets, it exhibited excellent generalization ability and strong
 491 robustness, which is instrumental in photovoltaic power generation prediction.

492 The space complexities of CloudRVE and ten alternative methods are summarized and compared in Table 8. It can
 493 be seen from the table that CloudRVE had a spatial complexity of 105.17 Mb, which is in line with the spatial
 494 complexity of Swin-T and Eff-Swin-T, and far less than the spatial complexity of ViT-L. The spatial complexity of
 495 CloudRVE exceeded that of RepVGG by three times, achieving the best ground cloud image classification
 496 performance. Thus, CloudRVE achieved excellent ground cloud image classification performance at the expense of
 497 higher spatial complexity.

498 **Table 8.** Space complexity of the proposed and ten alternative methods.

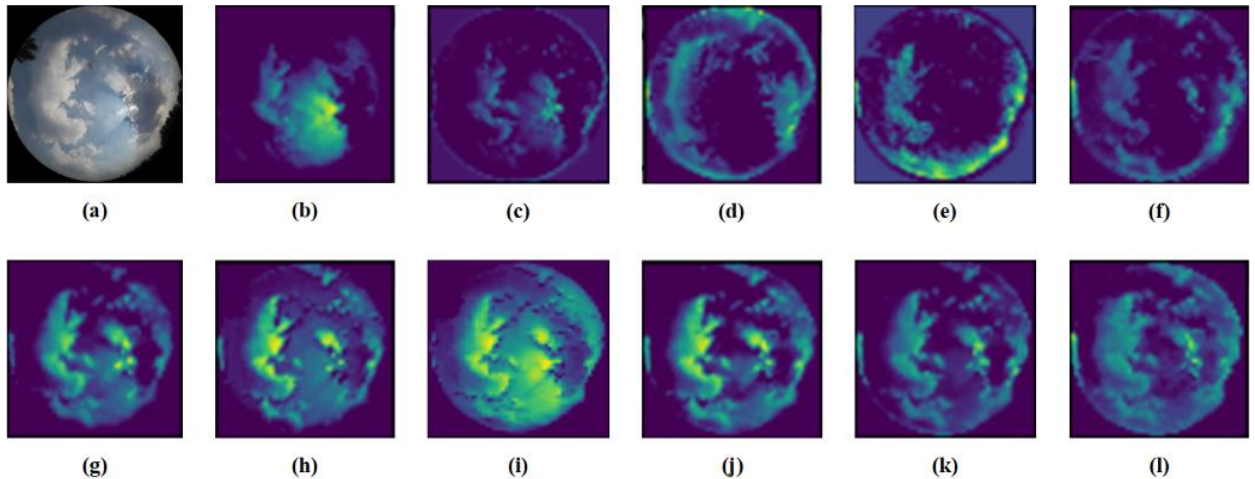
Method	Space complexity (Mb)
VGG-16	512.28
ResNet-50	90.03
ShuffleNet	4.93
CloudNet	153.36
CloudA	87.57
EfficientNet	15.61
ViT-L	327.37
Swin-T	105.28
RepVGG	30.10
Eff-Swin-T	105.24
CloudRVE	105.17

499 In order to provide a more intuitive display of the advantages of CloudRVE over other advanced methods, we
 500 extracted the features of the intermediate layers of different methods to generate the ground cloud feature maps for
 501 the building foundation, demonstrating the strong feature extraction capabilities of CloudRVE and proving its
 502 superiority, as shown in Figures 17 and 18. Feature extraction was achieved by generating rough feature maps through
 503 network training with parameter weights to highlight the important regions of predicted images. The light colored
 504 regions represent the important features, while the dark colored regions represent the sky or unsuccessfully extracted
 505 features. Figure 17(b-i) shows the feature maps of different ground cloud classification methods based on MGCD
 506 dataset to demonstrate the CloudRVE capability to extract more extensive and comprehensive cloud features and
 507 suppress the black regions and sunlight, further illustrating the best feature localization and extraction capability of
 508 CloudRVE. Figure 18(b-i) shows the feature maps of different ground cloud classification methods based on GRSCD
 509 dataset to demonstrate that the cloud feature extracted by CloudRVE covers the effective area in Figure 18(a) with
 510 the best coverage and the best suppression of the sunlight, further proving that CloudRVE has the best feature
 511 localization and extraction capabilities.



512

513 **Figure 17: Feature extraction of different methods based on MGCD, (a) Original (Liu et al., 2020a); (b)VGG-16; (c) ResNet-**
 514 **50; (d) ShuffleNet; (e) CloudNet; (f) CloudA; (g) EfficientNet; (h) ViT-L; (i) Swin-T; (j) RepVGG; (k) Eff-Swin-T; (l)**
 515 **CloudRVE**



516

517 **Figure 18: Feature extraction of different methods based on GRSCD: (a) Original (Liu et al., 2020b); (b)VGG-16;**
 518 **(c)ResNet-50; (d) ShuffleNet; (e) CloudNet; (f) CloudA; (g) EfficientNet; (h) ViT-L; (i) Swin-T; (j) RepVGG; (k) Eff-Swin-**
 519 **T; (l) CloudRVE**

520 5. Conclusion

521 This study proposed a new classification method called CloudRVE for ground-based cloud images based on the
 522 improved RepVGG network. In particular, its training stage structure was improved, the residual structure was
 523 broadened, and 1×1 convolutional layer branches were added to each block, extending the gradient information of the
 524 topology structure and enhancing the network ability to represent boundary features of cloud images. In addition, the
 525 NECA module was embedded after multi-branch fusion to learn the feature relationship between sequences, improve

526 the network cross-channel interaction ability, and extract the best global features of cloud images. We validated the
527 excellent performance of the proposed method on MGCD and GRSCD ground-based cloud image datasets, achieving
528 the classification accuracy values of 98.15 and 98.07%, respectively, which outperformed ten other advanced methods.
529 In addition, the MGCD and GRSCD ground-based cloud image datasets contain 7 types of cloud categories, which is
530 more than the ground-based cloud image datasets used in other papers. This further demonstrates the excellent
531 performance of the proposed method. The particular contributions of this paper were summarized in Section 1.
532 However, this study shares some limitations with other methods of classifying ground-based cloud images via
533 convolutional neural networks, which have reached a bottleneck due to continuous expansion of the capacity of
534 ground-based cloud image datasets. A lucrative alternative is Transformer, which got a high reputation of a powerful
535 deep neural network for processing sequences but has received little attention in ground-based cloud image
536 classification. On the other hand, cloud classification is only based on ground-based cloud image features, while many
537 physical features, such as height, thickness, etc., may be also used. Our follow-up study envisages combining CNN
538 and Transformer models and using cloud height, cloud thickness, and other parameters in ground-based cloud image
539 classification to improve the model's performance.

540

541 **Author Contributions.** LH performed the experiments and wrote the paper. CS, KZ, and HX analyzed the data and
542 designed the experiments. CS conceived the method and reviewed the paper. XL, ZS, and XZ reviewed the paper and
543 gave constructive suggestions.

544 **Financial support.** This research was funded by the National Science Foundation of China (NSFC) under Grant No.
545 62076093 and No. 62206095 and through the Fundamental Research Funds for the Central Universities of China
546 under Grant No. 2022MS078 and No. 2020MS099.

547

548 **Data Availability Statement.** The MGCD dataset was accessed from [https://github.com/shuangliutjnu/Multimodal-](https://github.com/shuangliutjnu/Multimodal-Ground-based-Cloud-Database)
549 [Ground-based-Cloud-Database](https://github.com/shuangliutjnu/Multimodal-Ground-based-Cloud-Database). The GRSCD dataset was accessed from [https://github.com/shuangliutjnu/TJNU-](https://github.com/shuangliutjnu/TJNU-Ground-based-Remote-Sensing-Cloud-Database)
550 [Ground-based-Remote-Sensing-Cloud-Database](https://github.com/shuangliutjnu/TJNU-Ground-based-Remote-Sensing-Cloud-Database).

551

552 **Acknowledgments.** We would like to thank Professor Liu Shuang of Tianjin Normal University for providing the
553 support of ground-based cloud image classification datasets and Student Meng Ru-oxuan from Guangxi Normal
554 University for her contribution to this paper.

555

556 **Declaration of Competing Interests.** The authors declare that they have no conflict of interest.

557 **References**

- 558 Alonso-Montesinos, J., Martinez-Durban, M., del Sagrado, J., del Aguila, I. M., and Batlles, F. J.: The application of
559 Bayesian network classifiers to cloud classification in satellite images, *Renew. Energy*, 97, 155–161,
560 <https://doi.org/10.1016/j.renene.2016.05.066>, 2016.
- 561 Calbó, J. and Sabburg, J.: Feature Extraction from Whole-Sky Ground-Based Images for Cloud-Type Recognition, *J.*
562 *Atmospheric Ocean. Technol.*, 25, 3–14, <https://doi.org/10.1175/2007JTECHA959.1>, 2008.
- 563 Cazorla, A., Olmo, F. J., and Alados-Arboledas, L.: Development of a sky imager for cloud cover assessment, *JOSA*
564 *A*, 25, 29–39, <https://doi.org/10.1364/JOSAA.25.000029>, 2008.
- 565 Ding, X., Zhang, X., Ma, N., Han, J., Ding, G., and Sun, J.: RepVGG: Making VGG-style ConvNets Great Again, in:
566 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021 IEEE/CVF Conference on
567 Computer Vision and Pattern Recognition (CVPR), 13728–13737, <https://doi.org/10.1109/CVPR46437.2021.01352>,
568 2021.
- 569 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M.,
570 Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N.: An Image is Worth 16x16 Words: Transformers for Image
571 Recognition at Scale, *International Conference on Learning Representations*, 1909, 2022.
- 572 Fabel, Y., Nouri, B., Wilbert, S., Blum, N., Triebel, R., Hasenbalg, M., Kuhn, P., Zarzalejo, L. F., and Pitz-Paal, R.:
573 Applying self-supervised learning for semantic cloud segmentation of all-sky images, *Atmospheric Meas. Tech.*, 15,
574 797–809, <https://doi.org/10.5194/amt-15-797-2022>, 2022.
- 575 Goren, T., Rosenfeld, D., Sourdeval, O., and Quaas, J.: Satellite Observations of Precipitating Marine Stratocumulus
576 Show Greater Cloud Fraction for Decoupled Clouds in Comparison to Coupled Clouds, *Geophys. Res. Lett.*, 45,
577 5126–5134, <https://doi.org/10.1029/2018GL078122>, 2018.
- 578 Gorodetskaya, I. V., Kneifel, S., Maahn, M., Van Tricht, K., Thiery, W., Schween, J. H., Mangold, A., Crewell, S.,
579 and Van Lipzig, N. P. M.: Cloud and precipitation properties from ground-based remote-sensing instruments in East
580 Antarctica, *The Cryosphere*, 9, 285–304, <https://doi.org/10.5194/tc-9-285-2015>, 2015.
- 581 Gyasi, E. K. and Swarnalatha, P.: Cloud-MobiNet: An Abridged Mobile-Net Convolutional Neural Network Model
582 for Ground-Based Cloud Classification, *Atmosphere*, 14, 280, <https://doi.org/10.3390/atmos14020280>, 2023.
- 583 He, K., Zhang, X., Ren, S., and Sun, J.: Deep Residual Learning for Image Recognition, in: 2016 IEEE Conference
584 on Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on Computer Vision and Pattern
585 Recognition (CVPR), 770–778, <https://doi.org/10.1109/CVPR.2016.90>, 2016.
- 586 He, T., Zhang, Z., Zhang, H., Zhang, Z., Xie, J., and Li, M.: Bag of Tricks for Image Classification with Convolutional
587 Neural Networks, 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Los Alamitos,
588 558–567, <https://doi.org/10.1109/CVPR.2019.00065>, 2019.

589 Heinle, A., Macke, A., and Srivastav, A.: Automatic cloud classification of whole sky images, *Atmospheric Meas.*
590 *Tech.*, 3, 557–567, <https://doi.org/10.5194/amt-3-557-2010>, 2010.

591 Hu, J., Shen, L., and Sun, G.: Squeeze-and-Excitation Networks, in: 2018 IEEE/CVF Conference on Computer Vision
592 and Pattern Recognition, 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 7132–7141,
593 <https://doi.org/10.1109/CVPR.2018.00745>, 2018.

594 Ioffe, S. and Szegedy, C.: Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate
595 Shift, in: *International Conference on Machine Learning*, Vol 37, San Diego, 448–456, 2015.

596 Kalisch, J. and Macke, A.: Estimation of the total cloud cover with high temporal resolution and parametrization of
597 short-term fluctuations of sea surface insolation, *Meteorol. Z.*, 603–611, <https://doi.org/10.1127/0941->
598 [2948/2008/0321](https://doi.org/10.1127/0941-2948/2008/0321), 2008.

599 Kingma, D. P. and Ba, J.: Adam: A Method for Stochastic Optimization, 3rd International Conference for Learning
600 Representations, San Diego, arXiv:1412.6980 [cs], <https://doi.org/10.48550/arXiv.1412.6980>, 2015.

601 Krizhevsky, A., Sutskever, I., and Hinton, G. E.: ImageNet classification with deep convolutional neural networks,
602 *Commun. ACM*, 60, 84–90, <https://doi.org/10.1145/3065386>, 2017.

603 Li, X., Qiu, B., Cao, G., Wu, C., and Zhang, L.: A Novel Method for Ground-Based Cloud Image Classification Using
604 Transformer, *Remote Sens.*, 14, 3978, <https://doi.org/10.3390/rs14163978>, 2022.

605 Li, Z., Kong, H., and Wong, C.-S.: Neural Network-Based Identification of Cloud Types from Ground-Based Images
606 of Cloud Layers, *Appl. Sci.*, 13, 4470, <https://doi.org/10.3390/app13074470>, 2023.

607 Lin, F., Zhang, Y., and Wang, J.: Recent advances in intra-hour solar forecasting: A review of ground-based sky
608 image methods, *Int. J. Forecast.*, 39, 244–265, <https://doi.org/10.1016/j.ijforecast.2021.11.002>, 2023.

609 Liu, S., Li, M., Zhang, Z., Cao, X., and Durrani, T. S.: Ground-Based Cloud Classification Using Task-Based Graph
610 Convolutional Network, *Geophys. Res. Lett.*, 47, e2020GL087338, <https://doi.org/10.1029/2020GL087338>, 2020a.

611 Liu, S., Li, M., Zhang, Z., Xiao, B., and Durrani, T. S.: Multi-Evidence and Multi-Modal Fusion Network for Ground-
612 Based Cloud Recognition, *Remote Sens.*, 12, 464, <https://doi.org/10.3390/rs12030464>, 2020b.

613 Liu, S., Duan, L., Zhang, Z., Cao, X., and Durrani, T. S.: Multimodal Ground-Based Remote Sensing Cloud
614 Classification via Learning Heterogeneous Deep Features, *IEEE Trans. Geosci. Remote Sens.*, 58, 7790–7800,
615 <https://doi.org/10.1109/TGRS.2020.2984265>, 2020c.

616 Liu, S., Duan, L., Zhang, Z., Cao, X., and Durrani, T. S.: Ground-Based Remote Sensing Cloud Classification via
617 Context Graph Attention Network, *IEEE Trans. Geosci. Remote Sens.*, 60, 1–11,
618 <https://doi.org/10.1109/TGRS.2021.3063255>, 2022.

619 Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B.: Swin Transformer: Hierarchical Vision
620 Transformer using Shifted Windows, in: 2021 IEEE/CVF International Conference on Computer Vision (ICCV),
621 2021 IEEE/CVF International Conference on Computer Vision (ICCV), Montreal, QC, Canada, 9992–10002,
622 <https://doi.org/10.1109/ICCV48922.2021.00986>, 2021.

623 Long, C., Li, X., Jing, Y., and Shen, H.: Bishift Networks for Thick Cloud Removal with Multitemporal Remote
624 Sensing Images, *Int. J. Intell. Syst.*, 2023, e9953198, <https://doi.org/10.1155/2023/9953198>, 2023.

625 Long, C. N., Sabburg, J. M., Calbó, J., and Pagès, D.: Retrieving Cloud Characteristics from Ground-Based Daytime
626 Color All-Sky Images, *J. Atmospheric Ocean. Technol.*, 23, 633–652, <https://doi.org/10.1175/JTECH1875.1>, 2006.

627 Meng, Q., Zhao, S., Huang, Z., and Zhou, F.: MagFace: A Universal Representation for Face Recognition and Quality
628 Assessment, in: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021 IEEE/CVF
629 Conference on Computer Vision and Pattern Recognition (CVPR), 14220–14229,
630 <https://doi.org/10.1109/CVPR46437.2021.01400>, 2021.

631 Norris, J. R., Allen, R. J., Evan, A. T., Zelinka, M. D., O’Dell, C. W., and Klein, S. A.: Evidence for climate change
632 in the satellite cloud record, *Nature*, 536, 72–75, <https://doi.org/10.1038/nature18273>, 2016.

633 Nouri, B., Kuhn, P., Wilbert, S., Hanrieder, N., Prah, C., Zarzalejo, L., Kazantzidis, A., Blanc, P., and Pitz-Paal, R.:
634 Cloud height and tracking accuracy of three all sky imager systems for individual clouds, *Sol. Energy*, 177, 213–228,
635 <https://doi.org/10.1016/j.solener.2018.10.079>, 2019.

636 Pfister, G., McKenzie, R. L., Liley, J. B., Thomas, A., Forgan, B. W., and Long, C. N.: Cloud Coverage Based on
637 All-Sky Imaging and Its Impact on Surface Solar Irradiance, *J. Appl. Meteorol. Climatol.*, 42, 1421–1434,
638 [https://doi.org/10.1175/1520-0450\(2003\)042<1421:CCBOAI>2.0.CO;2](https://doi.org/10.1175/1520-0450(2003)042<1421:CCBOAI>2.0.CO;2), 2003.

639 Qu, Y., Xu, J., Sun, Y., and Liu, D.: A temporal distributed hybrid deep learning model for day-ahead distributed PV
640 power forecasting, *Appl. Energy*, 304, 117704, <https://doi.org/10.1016/j.apenergy.2021.117704>, 2021.

641 Sarukkai, V., Jain, A., UzKent, B., and Ermon, S.: Cloud Removal in Satellite Images Using Spatiotemporal
642 Generative Networks, in: 2020 IEEE Winter Conference on Applications of Computer Vision (WACV), 2020 IEEE
643 Winter Conference on Applications of Computer Vision (WACV), 1785–1794,
644 <https://doi.org/10.1109/WACV45572.2020.9093564>, 2020.

645 Shi, C., Wang, C., Wang, Y., and Xiao, B.: Deep Convolutional Activations-Based Features for Ground-Based Cloud
646 Classification, *IEEE Geosci. Remote Sens. Lett.*, 14, 816–820, <https://doi.org/10.1109/LGRS.2017.2681658>, 2017.

647 Shi, C., Zhou, Y., Qiu, B., He, J., Ding, M., and Wei, S.: Diurnal and nocturnal cloud segmentation of all-sky imager
648 (ASI) images using enhancement fully convolutional networks, *Atmospheric Meas. Tech.*, 12, 4713–4724,
649 <https://doi.org/10.5194/amt-12-4713-2019>, 2019.

650 Simonyan, K. and Zisserman, A.: Very Deep Convolutional Networks for Large-Scale Image Recognition,
651 <https://doi.org/10.48550/arXiv.1409.1556>, 10 April 2015.

652 Singh, M. and Glennen, M.: Automated ground-based cloud recognition, *Pattern Anal. Appl.*, 8, 258–271,
653 <https://doi.org/10.1007/s10044-005-0007-5>, 2005.

654 Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A.:
655 Going deeper with convolutions, in: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR),

656 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 1–9,
657 <https://doi.org/10.1109/CVPR.2015.7298594>, 2015.

658 Tan, M. and Le, Q. V.: EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks, in: International
659 Conference on Machine Learning, Vol 97, San Diego, 2019.

660 Tang, Y., Yang, P., Zhou, Z., Pan, D., Chen, J., and Zhao, X.: Improving cloud type classification of ground-based
661 images using region covariance descriptors, *Atmospheric Meas. Tech.*, 14, 737–747, [https://doi.org/10.5194/amt-14-](https://doi.org/10.5194/amt-14-737-2021)
662 [737-2021](https://doi.org/10.5194/amt-14-737-2021), 2021.

663 Taravat, A., Del Frate, F., Cornaro, C., and Vergari, S.: Neural Networks and Support Vector Machine Algorithms
664 for Automatic Cloud Classification of Whole-Sky Ground-Based Images, *IEEE Geosci. Remote Sens. Lett.*, 12, 666–
665 670, <https://doi.org/10.1109/LGRS.2014.2356616>, 2015.

666 Wang, M., Zhou, S., Yang, Z., and Liu, Z.: CloudA: A Ground-Based Cloud Classification Method with a
667 Convolutional Neural Network, *J. Atmospheric Ocean. Technol.*, 37, 1661–1668, [https://doi.org/10.1175/JTECH-D-](https://doi.org/10.1175/JTECH-D-19-0189.1)
668 [19-0189.1](https://doi.org/10.1175/JTECH-D-19-0189.1), 2020.

669 Wang, M., Zhuang, Z., Wang, K., Zhou, S., Zhou, S., and Liu, Z.: Intelligent classification of ground-based visible
670 cloud images using a transfer convolutional neural network and fine-tuning, *Opt. Express*, 29, 41176–41190,
671 <https://doi.org/10.1364/OE.442455>, 2021.

672 Wu, X., Zhan, C., Lai, Y.-K., Cheng, M.-M., and Yang, J.: A Large-Scale Benchmark Dataset for Insect Pest
673 Recognition, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019 IEEE/CVF
674 Conference on Computer Vision and Pattern Recognition (CVPR), Conference Location: Long Beach, CA, USA,
675 8779–8788, <https://doi.org/10.1109/CVPR.2019.00899>, 2019.

676 Ye, L., Cao, Z., and Xiao, Y.: DeepCloud: Ground-Based Cloud Image Categorization Using Deep Convolutional
677 Features, *IEEE Trans. Geosci. Remote Sens.*, 55, 5729–5740, <https://doi.org/10.1109/TGRS.2017.2712809>, 2017.

678 Yu, A., Tang, M., Li, G., Hou, B., Xuan, Z., Zhu, B., and Chen, T.: A Novel Robust Classification Method for Ground-
679 Based Clouds, *Atmosphere*, 12, 999, <https://doi.org/10.3390/atmos12080999>, 2021.

680 Zhang, J., Liu, P., Zhang, F., and Song, Q.: CloudNet: Ground-Based Cloud Classification With Deep Convolutional
681 Neural Network, *Geophys. Res. Lett.*, 45, 8665–8672, <https://doi.org/10.1029/2018GL077787>, 2018a.

682 Zhang, X., Zhou, X., Lin, M., and Sun, R.: ShuffleNet: An Extremely Efficient Convolutional Neural Network for
683 Mobile Devices, in: 2018 Ieee/Cvf Conference on Computer Vision and Pattern Recognition (cvpr), New York, 6848–
684 6856, <https://doi.org/10.1109/CVPR.2018.00716>, 2018b.

685 Zhang, Y., Liu, H., and Hu, Q.: TransFuse: Fusing Transformers and CNNs for Medical Image Segmentation, in:
686 Medical Image Computing and Computer Assisted Intervention – MICCAI 2021, Cham, 14–24,
687 https://doi.org/10.1007/978-3-030-87193-2_2, 2021.

688 Zhao, Z., Xu, G., Qi, Y., Liu, N., and Zhang, T.: Multi-patch deep features for power line insulator status classification
689 from aerial images, in: 2016 International Joint Conference on Neural Networks (IJCNN), 2016 International Joint
690 Conference on Neural Networks (IJCNN), 3187–3194, <https://doi.org/10.1109/IJCNN.2016.7727606>, 2016.

691 Zheng, Y., Rosenfeld, D., Zhu, Y., and Li, Z.: Satellite-Based Estimation of Cloud Top Radiative Cooling Rate for
692 Marine Stratocumulus, *Geophys. Res. Lett.*, 46, 4485–4494, <https://doi.org/10.1029/2019GL082094>, 2019.

693 Zhong, B., Chen, W., Wu, S., Hu, L., Luo, X., and Liu, Q.: A Cloud Detection Method Based on Relationship Between
694 Objects of Cloud and Cloud-Shadow for Chinese Moderate to High Resolution Satellite Imagery, *IEEE J. Sel. Top.*
695 *Appl. Earth Obs. Remote Sens.*, 10, 4898–4908, <https://doi.org/10.1109/JSTARS.2017.2734912>, 2017.

696 Zhu, W., Chen, T., Hou, B., Bian, C., Yu, A., Chen, L., Tang, M., and Zhu, Y.: Classification of Ground-Based Cloud
697 Images by Improved Combined Convolutional Network, *Appl. Sci.*, 12, 1570, <https://doi.org/10.3390/app12031570>,
698 2022.

699 Zhuo, W., Cao, Z., and Xiao, Y.: Cloud Classification of Ground-Based Images Using Texture–Structure Features, *J.*
700 *Atmospheric Ocean. Technol.*, 31, 79–92, <https://doi.org/10.1175/JTECH-D-13-00048.1>, 2014.

701