

1       **Retrieval of Cloud Fraction using Machine Learning Algorithms**  
2                                   **based on FY4A AGRI observations.**

3                                   Jinyi Xia<sup>1</sup>       Li Guan<sup>1</sup>

4       <sup>1</sup>China Meteorological Administration Aerosol-Cloud and Precipitation Key  
5       Laboratory, Nanjing University of Information Science and Technology, Nanjing  
6       210044, China

7       Correspondence to: Li Guan   [liguan@nuist.edu.cn](mailto:liguan@nuist.edu.cn)

8       **Abstract**

9       Cloud fraction as a vital component of meteorological satellite products plays an  
10       essential role in environmental monitoring, disaster detection, climate analysis, and  
11       other research areas. Random Forest(RF) and Multilayer Perceptron(MLP) algorithms  
12       were used in this paper to retrieve the cloud fraction of AGRI (Advanced  
13       Geosynchronous Radiation Imager) onboard FY-4A satellite based on its full-disc level-  
14       1 radiance observation. Corrections has been made subsequently to the retrieved cloud  
15       fraction in areas where solar glint occurs using a correction curve fitted with sun-glint  
16       angle as weight. The algorithm includes two steps: the cloud detection is conducted  
17       firstly for each AGRI field of view to identify whether it is clear sky, partly cloudy or  
18       overcast within the observation field. Then the cloud fraction is retrieved for the scene  
19       identified as partly cloudy. The 2B-CLDCLASS-LIDAR cloud fraction product from  
20       Cloudsat& CALIPSO active remote sensing satellite is employed as the truth to assess  
21       the accuracy of the retrieval algorithm. Comparison with the operational AGRI level 2  
22       cloud fraction product is also conducted at the same time. The results indicate that both  
23       the Random Forest (RF) and Multi-Layer Perceptron (MLP) cloud detection models  
24       achieved high accuracy, surpassing that of operational products. However, both  
25       algorithms demonstrated weaker discrimination capabilities for partly cloudy  
26       conditions compared to clear sky and overcast situations. Specifically, they tended to  
27       misclassify fields of view with low cloud fractions (e.g., cloud fraction = 0.16) as clear

28 sky and those with higher cloud fractions (e.g., cloud fraction = 0.83) as overcast.  
29 Between the two models, RF exhibited higher overall accuracy. Both RF and MLP  
30 models performed well in cloud fraction retrieval, showing lower mean error (ME),  
31 mean absolute error (MAE), and root mean square error (RMSE) compared to  
32 operational products. The ME for both RF and MLP cloud fraction retrieval models was  
33 close to zero, while RF had slightly lower MAE and RMSE than MLP. During daytime,  
34 the high reflectance in sun-glint areas led to larger retrieval errors for both RF and MLP  
35 algorithms. However, after correction, the retrieval accuracy in these regions improved  
36 significantly. At night, the absence of visible light observations from the AGRI  
37 instrument resulted in lower classification accuracy compared to daytime, leading to  
38 higher cloud fraction retrieval errors during nighttime.

39 **Key words:** Cloud detection; cloud fraction retrieval; FY-4A AGRI; CloudSat &  
40 CALIPSO; machine learning; deep learning.

## 41 **Introduction**

42 Clouds occupy a significant proportion within satellite remote sensing data  
43 acquired for Earth observation. According to the statistics from the International  
44 Satellite Cloud Climatology Project (ISCCP), the annual average global cloud coverage  
45 within satellite remote sensing data is around 66% with even higher cloud coverage in  
46 specific regions (such as the tropics) (Zhang, et al., 2004). The impact of clouds on the  
47 radiation balance of the Earth's atmospheric system is influenced by the optical  
48 properties of clouds. Cloud detection, as a vital component of remote sensing image  
49 data processing, is considered a critical step for the subsequent identification, analysis,  
50 and interpretation of remote sensing images. Therefore, accurately determining cloud  
51 coverage is essential in various research domains, such as environmental monitoring,  
52 disaster surveillance and climate analysis.

53 Fengyun-4A (FY-4A) is a comprehensive atmospheric observation satellite

54 launched by China in 2016. The uploaded AGRI (Advanced Geosynchronous Radiation  
55 Imager) has 14 channels and captures full-disk observation every 15 minutes. In  
56 addition to observing clouds, water vapor, vegetation and the Earth's surface, it also  
57 possesses the capability to capture aerosols and snow. Moreover, it can clearly  
58 distinguish different phases and particle size of clouds and obtain high- to mid-level  
59 water vapor content. It is particularly suitable for cloud detection due to its  
60 simultaneous use of visible, near-infrared, and long-wave infrared channels for  
61 observation with 4km spatial resolution.

62 Numerous cloud detection algorithms have been provided based on observations  
63 from satellite-borne imagers. The threshold method has been widely employed by  
64 researchers, including the early ISCCP (International Satellite Cloud Climatology  
65 Project) method (Rossow, 1993) and the proposed threshold methods based on different  
66 spectral features or underlying surfaces (Kegelmeyer,1994; Solvsteen,1995; Baum and  
67 Trepte,1996). However, there is a significant subjectivity in selection of thresholds  
68 whether it is the single and fixed threshold in the early days, multiple thresholds,  
69 dynamic thresholds, or adaptive thresholds. The selection of thresholds is influenced  
70 by season and climate. Surface reflectance varies significantly between different  
71 seasons, such as increased reflectance from snow in winter and vegetation flourishing  
72 in summer affecting reflectance. As a result, changes in surface features during different  
73 seasons lead to variations in the distribution of grayscale values in images, requiring  
74 adjustments to thresholds based on seasonal characteristics. Climate conditions like  
75 cloud cover, atmospheric humidity, etc., impact the distinguishability of clouds and  
76 other features. For instance, in humid or cloudy climates, the reflectance of the surface  
77 and clouds may be similar, necessitating stricter thresholds for differentiation.  
78 Therefore, climate conditions also influence threshold selection.

79 The other category of cloud detection algorithms is based on statistical probability  
80 theory. For example the principal component discriminant analysis and quadratic

81 discriminant analysis methods were used to SEVIRI (Spinning Enhanced Visible and  
82 Infrared Imager) cloud detection (Amato et al., 2008). The cloud detection algorithm  
83 for Thermal Infrared (TIR) sensor was based on the Bayesian theory of total probability  
84 (Merchant et al., 2010) and the naive Bayes algorithm for AGRI (Yan , et al., 2022).  
85 The unsupervised clustering cloud detection algorithms for MERIS (Medium  
86 Resolution Imaging Spectrometer) (GomezChova , et al., 2007) and the fuzzy C-means  
87 clustering algorithms for MODIS (Pan, et al., 2009) all have achieved high accuracy in  
88 cloud detection.

89 More and more machine learning algorithms are being utilized by researchers in  
90 cloud detection studies with the development of machine learning. For instance, the  
91 probabilistic neural networks, especially radial basis function networks was used for  
92 AVHRR cloud detection (Zhang, et al., 2001). The utilization of convolutional neural  
93 network methods (Chai, et al., 2024) offers important perspectives for cloud detection  
94 research.

95 Currently, there is limited research literature on cloud detection and cloud fraction  
96 retrieval algorithms for FY-4A/4B AGRI. The operational cloud fraction product of FY-  
97 4A AGRI utilized a threshold method with 4 km spatial resolution. Differences in  
98 climatic and environmental factors lead to varying albedo and brightness temperature  
99 observations for the instrument at different times and locations. Therefore, the choice  
100 of thresholds is easily influenced by factors such as season, latitude and land surface  
101 type (Gao and Jing, 2019). Using multiple sets of thresholds for discrimination would  
102 significantly slow down the cloud detection process. Moreover, most algorithms focus  
103 solely on cloud detection, which classified the observed scenes into cloud or clear-sky  
104 without providing the specific cloud fraction information for the scenes. The use of  
105 active remote sensing instruments carried by Cloudsat & Calypso is not influenced by  
106 thresholds when retrieving cloud fraction, enabling a more accurate cloud fraction  
107 retrieval. However, due to Cloudsat & Calypso being polar-orbiting satellites, the cloud

108 fraction over the full disk cannot be obtained. Utilizing the Cloudsat & Calypso Level  
109 2 product 2B-CLDCLASS-LIDAR as the reference truth, a random forest model trained  
110 based on FY4A AGRI full disk radiation data can address the shortcomings of threshold  
111 methods and achieve a high accuracy of cloud fraction over the full disk.

112 In summary, this paper established cloud detection and cloud fraction retrieval  
113 models using a Multi-Layer Perceptron (MLP) and Random Forest (RF), based on FY-  
114 4A AGRI full-disk level 1 observed radiance data. The cloud fraction from the CloudSat  
115 & CALIPSO level 2 product 2B-CLDCLASS-LIDAR was used as the label. The results  
116 were compared with the 2B-CLDCLASS-LIDAR product and the official AGRI  
117 operational products for validation.

## 118 **1 Research Data and Preprocessing**

### 119 *1.1 FY-4A data*

120 FY-4A was successfully launched on December 11, 2016. Starting from May 25, 2017,  
121 FY-4A drifted to a position near the main business location of the Fengyun  
122 geostationary satellite at 104.7 degrees east longitude on the equator. Its successful  
123 launch marked the beginning of a new era for China's next-generation geostationary  
124 meteorological satellites as an advanced comprehensive atmospheric observation  
125 satellite. The Advanced Geosynchronous Radiation Imager (AGRI), one of the main  
126 payloads of the Fengyun-4 series geostationary meteorological satellites, can perform  
127 large-disk scans and rapid regional scans at a minute level. It has 14 observation  
128 channels in total with the main task of acquiring cloud images. The channel parameters  
129 and main uses of AGRI are detailed in Table 1  
130 (<https://www.nsmc.org.cn/nsmc/cn/instrument/AGRI.html>). The first six visible light  
131 channels have no values at night, meaning that channels with a central wavelength less  
132 than or equal to 2.225 $\mu$ m are unavailable during nighttime. FY-4A AGRI data was  
133 downloaded from the official website of the China national satellite meteorological

134 center (<http://satellite.nsmc.org.cn>), including level-1 full disk radiation observation  
 135 data preprocessed through quality control, geolocation and radiation calibration as well  
 136 as level-2 cloud fraction product (CFR). The spatial resolution of these data is all 4 km  
 137 at nadir and the temporal resolution is 15 minutes.

138 Table 1 FY-4A AGRI channel parameters

Channel Number	Band Range / $\mu\text{m}$	Central Wavelength / $\mu\text{m}$	Spatial resolution/km	Main Applications
1	0.45 ~ 0.49	0.47	1	clouds, dust, aerosols
2	0.55 ~ 0.75	0.65	0.5	clouds, sand dust, snow
3	0.75 ~ 0.90	0.825	1	vegetation
4	1.36 ~ 1.39	1.375	2	cirrus
5	1.58 ~ 1.64	1.61	2	clouds、 snow
6	2.10 ~ 2.35	2.225	2	cirrus、 aerosols
7	3.50 ~ 4.00	3.75H	2	fire point, the intense solar reflection signal
8	3.50 ~ 4.00	3.75L	4	low clouds, fog
9	5.80 ~ 6.70	6.25	4	upper-level water vapor
10	6.90 ~ 7.30	7.1	4	mid-level water vapor
11	8.00 ~ 9.00	8.5	4	subsurface water vapor
12	10.30 ~ 11.30	10.8	4	surface and cloud-top temperatures
e	11.50 ~ 12.50	12.0	4	surface and cloud-top temperatures
14	13.2 ~ 13.8	13.5	4	cloud-top height

139 **1.2 CloudSat & Calipso Cloud Product**

140 CALIPSO (Cloud-Aerosol Lidar and Infrared Pathfinder Satellite Observations)  
 141 is a satellite jointly launched by NASA and CNES (the French National Center for  
 142 Space Studies) in 2006. It is a member of the A-Train satellite observation system.  
 143 CALIPSO is equipped with three payloads, among which CALIOP (the Cloud and  
 144 Aerosol Lidar with Orthogonal Polarization) is a primary observational instrument.  
 145 Observing with dual wavelengths (532 nm and 1064 nm) CALIOP can provide high-  
 146 resolution vertical profiles of clouds and aerosols with 30 m vertical resolution. As the  
 147 first satellite designed to observe global cloud characteristics in a sun-synchronous orbit  
 148 CloudSat is also among NASA's A-Train series satellites. The CPR (Cloud Profile  
 149 Radar) installed on it operates at 94 GHz millimeter-wave and is capable of detecting  
 150 the vertical structure of clouds and providing vertical profiles of cloud parameters. The

151 scanning wavelengths of CPR and CALIOP are different. CALIOP is capable of  
152 observing the top of mid-to-high level clouds, whereas CPR can penetrate optically  
153 thick clouds. Combining the strengths of these two instruments enables the acquisition  
154 of precise and detailed information on cloud layers and cloud fraction.

155 The joint level 2 product 2B-CLDCLASS-LIDAR is mainly utilizing in this study.  
156 It provides the cloud fraction at different heights with horizontal resolution 2.5 km  
157 (along-track)  $\times$  1.4 km (cross-track) through combining the observations from CPR and  
158 CALIOP. Since the two instruments have different spatial domain such as vertical  
159 resolution, spatial resolution and spatial frequency, the spatial domain of the output  
160 products is defined in terms of the spatial grid of the CPR. In the algorithm, the cloud  
161 fraction is calculated using a weighted scheme based on the spatial probability of  
162 overlap between the radar and lidar observations. The calculation of the lidar cloud  
163 fraction within a radar footprint is represented by the equation 1 (Mace, G. G., et al,  
164 2007):

$$165 \quad C_l = \frac{\sum_{i=1}^{\# \text{ of lidar obs}} w_i \delta_i}{\sum_{i=1}^{\# \text{ of lidar obs}} w_i} \quad (1)$$

166 Where:

167  $C_l$  represents the lidar cloud fraction within a radar footprint.

168  $w_i$  is the spatial probability of overlap for a particular lidar observation.

169  $\delta_i$  indicates the lidar hydrometeor occurrence, where a value of 1 signifies the  
170 presence of hydrometeor and 0 indicates the absence.

171  $i$  counts the lidar profile in a specific radar observational domain.

172 This calculation considers the contributions of multiple lidar observations within  
173 a radar resolution volume to determine the cloud fraction within that volume. The  
174 CloudSat product manual (Wang, 2019) can be referred for more detailed information  
175 on 2B-CLDCLASS-LIDAR. The data used is available to download from the ICARE  
176 data and services center ([https://www.icare.univ-lille.fr/data-access/data-archive-  
177 access/](https://www.icare.univ-lille.fr/data-access/data-archive-access/)).

### 178 *1.3 Establishment of Training Data*

179 The crucial aspect of establishing a training data in machine learning algorithms  
180 is how to obtain the cloud fraction values (ground truth) as labels. The error in cloud  
181 fraction retrieved solely from passive remote sensing instruments is significant. Using  
182 active remote sensing data can provide more accurate cloud fraction information in the  
183 vertical direction. Therefore, the spatiotemporally matched 2B-CLDCLASS-LIDAR  
184 cloud fraction are utilized as output labels in this paper.

185 The FY-4A AGRI and 2B-CLDCLASS-LIDAR data with a spatial difference  
186 between fields of view within 1.5 km and a time difference within 15 minutes are  
187 spatiotemporal matched. To make the 2B-CLDCLASS-LIDAR cloud fraction data  
188 collocated within AGRI pixels more effective, at least two 2B-CLDCLASS-LIDAR  
189 pixels are required within each AGRI field of view. The cloud fraction average of these  
190 pixels is used as the cloud fraction for that AGRI pixel. However, the errors in the  
191 matched dataset are unavoidable. The AGRI scanning method operates from left to right  
192 and top to bottom. Each complete scan of the full disk takes 15 minutes and generates  
193 a dataset. It is impossible to determine the exact moment of a specific point within the  
194 full disk. This limits the time range for matching datasets to within 15 minutes.  
195 However, in areas with higher wind speeds, clouds can move a significant distance  
196 within that 15-minute window. Therefore, errors arising from timing issues cannot be  
197 avoided.

198 Cloud detection and cloud fraction label generation for 2B-CLDCLASS-LIDAR  
199 are as follows. There may be multiple layers of clouds in each field of view. If there is  
200 at least one layer cloud with cloud fraction of 1 in the 2B-CLDCLASS-LIDAR profile,  
201 then the scene is labeled as overcast with a cloud fraction of 1. If all layers in the profile  
202 are cloud-free, the scene is labeled as clear sky. The scene between the above two  
203 situations is labeled as partly cloudy and the cloud fraction is the average of cloud  
204 fractions at different layers.



205       The algorithm includes two steps: the cloud detection is conducted firstly for each  
206 AGRI field of view to identify whether it is clear sky, partly cloudy or overcast within  
207 the observation field. Then the cloud fraction is retrieved for the scene identified as  
208 partly cloudy. So the training data include dataset A used for cloud detection and dataset  
209 B for cloud fraction retrieval. The input variables in dataset A are the FY-4A AGRI  
210 level-1 radiative observations from 14 channels and the output variable is the  
211 temporally and spatially matched 2B-CLDCLASS-LIDAR cloud detection label. The  
212 output is categorized into three types: overcast, partly cloudy and clear sky with values  
213 1, 2 and 3 respectively. The cloud fraction product from 2B-CLDCLASS-LIDAR  
214 consists of discrete values: 0, 0.16, 0.33, 0.50, 0.66, 0.83, and 1. According to the result  
215 statistics, the cloud fractions of 2B-CLDCLASS-LIDAR pixels within the AGRI field  
216 of view are mostly the same. After averaging, the proportions of cloud fractions of [0.16,  
217 0.33, 0.5, 0.67, 0.83] are extremely high. Therefore, other cloud fraction situations with  
218 extremely small proportions can be ignored. Doing so can also better balance the  
219 training samples. Here, 0 indicates clear sky, values from 0 to 1 represent varying cloud  
220 fractions for partly cloudy conditions, and 1 signifies overcast. To ensure the balance  
221 and representativeness of the samples, the proportions of different cloud fraction  
222 samples in dataset A are set at 5:1:1:1:1:5. Regarding the samples for partly cloudy  
223 type in dataset A, the collocated 2B-CLDCLASS-LIDAR cloud fraction products serve  
224 as output labels for cloud fraction retrieval model B. The input of training dataset B  
225 remains the FY-4A AGRI level-1 radiative observations.

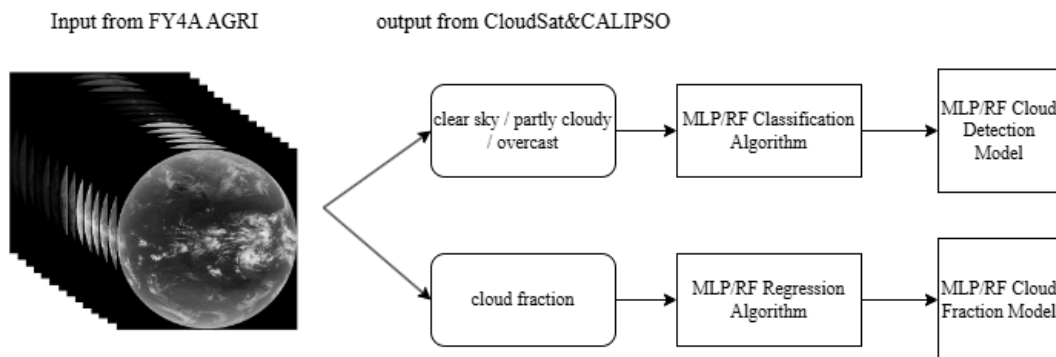
226       Due to the instrument's limited lifespan, only 2B-CLDCLASS-LIDAR data up to  
227 August 2019 can be obtained. The sample time range used in this paper is from August  
228 2018 to July 2019. Five days were randomly selected each month as daytime samples  
229 and five days as nighttime samples. A total of 120 days of time and space matched FY-  
230 4A AGRI full-disk observations and 2B-CLDCLASS-LIDAR data were used as  
231 training and testing samples. Among them, 80% of the data was used for training, and

232 20% was used for testing. The total number of daytime samples in dataset A is 91,073,  
 233 while dataset B contains 30,358 samples. The total number of nighttime samples in  
 234 dataset A is 95,493, and dataset B includes 31,831 samples.

235 Although the model was trained and tested using data from 2018 to 2019, to test  
 236 the universality of the algorithm, it was applied to real-time observations from FY-4A  
 237 and FY-4B AGRI in 2023.

## 238 2 Algorithms

239 Our preliminary experiments involved multiple algorithms, including LibSvm,  
 240 MLP, BP neural network, and Random Forest. These experiments highlighted that,  
 241 among the baselines, Random Forest and MLP achieved the highest overall accuracy.  
 242 For this reason, we selected them to perform additional experiments. Using RF and  
 243 MLP algorithms to train the model with the established sample set, the overall process  
 244 is shown in the Figure 1.



245 Observations of 14 channels for each pixel.

246 **Figure 1:** Method workflow. The input consists of 14 channel observation values  
 247 for each pixel from FY4A AGRI, and the ground truth labels or outputs are sourced  
 248 from the CloudSat&CALIPSO cloud fraction products. The cloud detection  
 249 classification model and the cloud fraction retrieval model are established separately.

### 250 2.1 Random Forest (RF)

251 This algorithm integrates multiple trees based on the Bagging idea of ensemble  
 252 learning, with the basic element being the decision tree (Breiman, 1999). When building

253 a decision tree,  $N$  sets of independent and dependent variables are randomly sampled  
254 with replacement from the original training samples to create a new training sample set;  
255  $m$  variables are randomly sampled without replacement from all independent variables,  
256 the dependent variable data is split into two parts using the selected variables, and the  
257 purity of the subsets is calculated for each split method. The variable utilized by the  
258 split method with the highest purity is used to partition the data, completing the decision  
259 at that node. This process of binary splitting continues to grow the decision tree until  
260 stopping criteria are met, completing the construction of a single decision tree. These  
261 steps are repeated  $N_{tree}$  times to build a random forest model consisting of  $N_{tree}$   
262 decision trees (Breiman, 2001). Random Forest adopts ensemble algorithms, with the  
263 advantage of high accuracy. It can handle both discrete and continuous data, without  
264 the need for normalization, making it more efficient compared to other algorithms.

## 265 **2.2 Multilayer Perceptron (MLP)**

266 This algorithm consists of a fully connected artificial neural network (Duda, et al.,  
267 2001). The classifier/regressor takes feature vectors or tensors as input. The input is  
268 mapped through multiple fully connected hidden layers containing hidden weights,  
269 which produce classifications/regressions at the output layer. A nonlinear activation  
270 function (such as sigmoid or rectified linear unit (ReLU)) is applied in each hidden  
271 layer to facilitate a nonlinear model. For classifiers, the output of the final hidden layer  
272 is combined and passed through a softmax function to generate class predictions. For  
273 the loss function, the cloud detection model is cross-entropy, and the cloud fraction  
274 model is MSE. The model's weights are trained in a supervised manner using  
275 backpropagation.

## 276 **2.3 Hyperparameters**

277 In this paper, a total of eight models were established, including daytime/nighttime

278 random forest classification/regression models and daytime/nighttime MLP  
279 classification/regression models. For the random forest, we first conducted experiments  
280 using the following Hyperparameters ranges: Trees: [200, 300, 400, 500, 600,700],  
281 minleaf: [1, 2, 5, 10], criterion: [Gini, entropy]. Ultimately, the best selections were: (1)  
282 Daytime RF classification model: Trees=500. (2) Nighttime RF classification model:  
283 Trees=600. (3) Daytime RF regression model: Trees=400. (4) Nighttime RF regression  
284 model: Trees=500. All four models have minleaf=1, criterion=gini.

285 For the MLP, experiments were conducted using the following hyperparameter  
286 ranges: Number of hidden layers: [2,3,4,5,6,7,8,9], Hidden layer size: [8,16,32,64,128],  
287 Epochs: [30,50,100], Solver hyperparameter: [lbfgs, sgd, adam]. The optimal  
288 parameters found are as follows: (1) MLP classification model for daytime: number of  
289 hidden layers = 5. (2) MLP classification model for nighttime: number of hidden layers  
290 = 5. (3) MLP regression model for daytime: number of hidden layers = 4. (4) MLP  
291 regression model for nighttime: number of hidden layers = 6. All four models have  
292 Hidden layer size = 64, Epochs = 50, solver = adam, BatchSize = 1500, Initial learning  
293 rate = 0.01, Learning rate schedule = piecewise, Factor for dropping the learning rate =  
294 0.1, Number of epochs for dropping the learning rate = 10.

### 295 **3 Results and Analysis**

296 To assess the accuracy and stability of the retrieval model, two types of validation  
297 methods are utilized. One way involves a direct comparison from images, qualitatively  
298 comparing the model's retrieval results and official cloud fraction products with AGRI  
299 observed cloud images. Another approach uses 2B-CLDCLASS-LIDAR as the ground  
300 truth and introduces five parameters for quantitative comparison: recall, false alarm rate  
301 (FAR), mean error (ME), mean absolute error (MAE), and root mean square error  
302 (RMSE). To evaluate the ability of operational products, RF, and MLP cloud detection  
303 models to distinguish overcast, partly cloudy, and clear sky, the recall is calculated using

304 the formula  $POD=TP/(TP+FN)$ , and the false alarm rate is calculated using the formula  
305  $FAR=FP/(TP+FP)$ . Taking the overcast scene as an example, TP represents the number  
306 of correctly identified overcast conditions, FN represents the number of overcast  
307 conditions misidentified as partly cloudy or clear sky, and FP represents the number of  
308 clear sky or partly cloudy conditions misidentified as overcast. When assessing the  
309 accuracy of operational products and cloud fraction models for the cloud fraction  
310 retrieval results of partly cloudy scenes, mean error (ME), mean absolute error (MAE),  
311 and root mean square error (RMSE) are used.

### 312 *3.1 Objective Analysis of Cloud Fraction Retrievals*

313 First, using the 2B-CLDCLASS-LIDAR cloud fraction product as the ground truth,  
314 we calculated the accuracy of the operational cloud detection products. The results are  
315 shown in columns 3-4 of Table 2. The samples used for this statistic are the same as  
316 those for testing the model below (20% of dataset A).

317 Based on the cloud detection model trained above, cloud detection experiments  
318 were conducted using the test samples from dataset A. The time-space matched 2B  
319 CLDCLASS-LIDAR cloud fraction product served as the ground truth to assess the  
320 accuracy of cloud detection. The results are shown in columns 5-8 of Table 2. During  
321 the day, the Random Forest model achieved an overall accuracy of 94.2%, while the  
322 MLP model had an overall accuracy of 93.7%. The Random Forest model exhibited  
323 slightly higher recall rates for clear skies, partly cloudy, and overcast conditions  
324 compared to the MLP model, and its FAR was lower as well. Both models performed  
325 poorly in recognizing partly cloudy conditions, as the models tended to classify true  
326 cloud fractions of 0.16 as clear skies and those of 0.83 as overcast. At night, the Random  
327 Forest model achieved an overall accuracy of 89.4%, while the MLP model had an  
328 accuracy of 88.7%. The Random Forest model had higher recall rates for clear skies  
329 and partly cloudy conditions compared to the MLP, while the recall rates for overcast

330 conditions were similar for both models. The FAR for the Random Forest model was  
 331 lower than that of the MLP. Overall, both the Random Forest and MLP models showed  
 332 higher classification accuracy for clear skies, partly cloudy, and overcast conditions  
 333 compared to operational products, with the Random Forest model performing better.

334 Table 2: Recall Rate , FAR of Operational Cloud Detection Products and multiple  
 335 models.

	Sky Classification	Daytime Product	Nighttime Product	Daytime RF	Nighttime RF	Daytime MLP	Nighttime MLP
POD	Clear Sky	0.6359	0.5781	0.964	0.919	0.959	0.905
	Partly cloudy	0.7174	0.7449	0.914	0.845	0.895	0.808
	Overcast	0.7736	0.7384	0.959	0.919	0.957	0.920
FAR	Clear Sky	0.1778	0.0934	0.047	0.102	0.064	0.131
	Partly cloudy	0.1819	0.2117	0.078	0.153	0.085	0.172
	Overcast	0.2499	0.2683	0.038	0.061	0.039	0.063

336  
 337 Based on the previous model's assessment of the field of view as partly cloudy, the  
 338 cloud fraction in this AGRI field of view is retrieved using the cloud fraction model  
 339 established earlier. For model evaluation, both the operational product and the 2B-  
 340 CLDCLASS-LIDAR cloud fraction product are classified as partly cloudy, with the  
 341 matched 2B-CLDCLASS-LIDAR cloud fraction product considered as the ground truth.  
 342 The average error, mean absolute error, and root mean square error for both daytime  
 343 and nighttime operational products and cloud fraction model retrieval (Table 3) are  
 344 calculated. It can be observed that the average errors of both models are close to 0  
 345 during both daytime and nighttime. The errors are smaller during the day than at night,  
 346 with the RF model exhibiting lower errors than the MLP model. In summary, the errors  
 347 of both models are smaller than those of the operational products, and the RF model  
 348 performs better in the cloud fraction retrieval task.

349 Table 3: Errors of Cloud Fraction

	Daytime Product	Nighttime Product	Daytime RF	Daytime MLP	Nighttime RF	Nighttime MLP
--	--------------------	----------------------	---------------	----------------	--------------	------------------

ME	0.1987	0.2121	0.0006	-0.0009	-0.0028	-0.0032
MAE	0.2279	0.2441	0.1011	0.1053	0.1221	0.1322
RMSE	0.2776	0.2938	0.1285	0.1332	0.1510	0.1623

350

351 Based on the experiments mentioned above, the performance of RF in cloud  
 352 detection and cloud fraction retrieval slightly outperforms that of MLP. Therefore,  
 353 subsequent experiments will utilize the RF algorithm.

### 354 *3.2 Cloud fraction correction in sun glint regions*

355 Sun glint refers to the bright areas created by the reflection of sunlight to the  
 356 sensors of observation systems (satellites or aircrafts). This phenomenon usually occurs  
 357 on extensive water surfaces, such as oceans lakes or rivers. This specular reflection of  
 358 sunlight will cause an increase in the reflected solar radiation received by onboard  
 359 sensors, manifested as an enhancement of white brightness in visible images. The  
 360 increase in visible channel observation albedo will affect various subsequent  
 361 applications of data, including cloud detection and cloud cover retrieval, etc.

362 The position of Sun glint area can be determined using the SunGlintAngle value  
 363 in the FY-4A GEO file. SunGlintAngle is defined as the angle between the satellite  
 364 observation direction or reflected radiation direction and the mirror reflection direction  
 365 on a calm surface (horizontal plane). It is generally accepted that the range of  
 366 SunGlintAngle  $< 15^\circ$  is easily affected by sun glint (Kay S, et al., 2009). The positions  
 367 of the SunGlintAngle contour lines at 5 and  $15^\circ$  are marked in Figure 1(a). It can be  
 368 observed that the edge of sun glint in Figure 1(a) essentially overlaps with the position  
 369 of SunGlintAngle =  $15^\circ$ . Thus, the region where SunGlintAngle  $< 15^\circ$  is defined as the  
 370 sun glint range in this paper and only the cloud fraction within this range will be

371 adjusted in the subsequent correction.

372 To correct the cloud fraction in the sun-glnt areas, we first identified the fields of  
373 view (FOVs) where sun-glnt occurred during FY-4A AGRI observations from August  
374 2018 to July 2019, totaling 1,476 FOVs. When matching the sample set of the sun glnt  
375 area, two issues need to be explained. 1) Cloud fraction is the average of cloud fractions  
376 of different layers: Among the matched pixels, only one-layer cloud and two-layer  
377 cloud appear. When there are two layers of cloud, there is always one layer with a cloud  
378 fraction of 1. According to the previous description, when there is one layer with a cloud  
379 fraction of 1, this pixel should be regarded as fully cloudy. 2) The average cloud fraction  
380 of at least two CloudSat & CALIPSO pixels is taken as the cloud fraction of the AGRI  
381 pixel: Due to the very small area of the sun glnt area, the matching is very difficult. If  
382 at least two CloudSat & CALIPSO pixels within an AGRI pixel are required, this will  
383 make the available sample size very small. Therefore, when making the sample set of  
384 the sun glnt area, only one CloudSat & CALIPSO pixel within an AGRI pixel is  
385 required. Due to the above two reasons, the true cloud fraction in the sample is a discrete  
386 value. Subsequently, a direct least squares fitting was conducted between the retrieved  
387 cloud fraction and the collocated 2B-CLDCLASS-LIDAR cloud fraction (ground truth).  
388 The scatter plot is illustrated in Figure 2(b), where x-axis is the 2B-CLDCLASS-  
389 LIDAR cloud fraction and y-axis is the model-retrieved cloud fraction. The blue line  
390 represents the curve (namely Eq.2) fitted by the least squares method between the  
391 retrievals and the truths. The thin dash line is the  $x=y$  line. It is evident that the retrieved  
392 cloud fraction is generally slightly overestimated.

393 Taking observations at 04:00 on 5 June 2019 as an example, Figure 2(c) presents  
394 the distribution of SunGlintAngle and the flight trajectory of the Cloudsat&Calypso  
395 satellite. White circles denote the sun glnt region with SunGlintAngle  $< 15^\circ$  and the  
396 white line represents the satellite flight track. As depicted in the figure, the majority of  
397 Cloudsat&Calypso flight trajectories do not pass through the central position of sun



398 glint area but instead traverse locations with larger SunGlintAngle values. The  
399 intensity of sun glint effect decreases with the increase of SunGlintAngle. This  
400 suggests that the true values for spatial and temporal matching mostly do not fall within  
401 the strongest sun glint region. From Figure 2(d), it can be seen that the impact of sun  
402 glint becomes stronger as SunGlintAngle decreasing, which results in a higher  
403 observation albedo. This further leads to the overestimated cloud fraction values in the  
404 retrieval. It is evident that the cloud fraction error is related to the value of  
405 SunGlintAngle and this influence is not considered in Eq. (2). Directly applying  
406 equation (2) to correct the cloud fraction retrievals would result in a too small correction  
407 intensity for the FOVs near the center of sun glint and an excessively large correction  
408 intensity for the FOVs in the Sun-glint edge region (even erroneous clear sky may  
409 appear). Considering this, a correction formula (3)-(4) using SunGlintAngle as weight  
410 is introduced, where  $W_i$  represents the angle weight for a certain pixel  $i$  in the sun glint  
411 region,  $n$  is the number of pixels within the SunGlintAngle  $< 15^\circ$  range,  $y_i$  is the initial  
412 model retrieval of cloud cover for the field of view  $i$  and  $x_i$  is the final corrected cloud  
413 fraction.

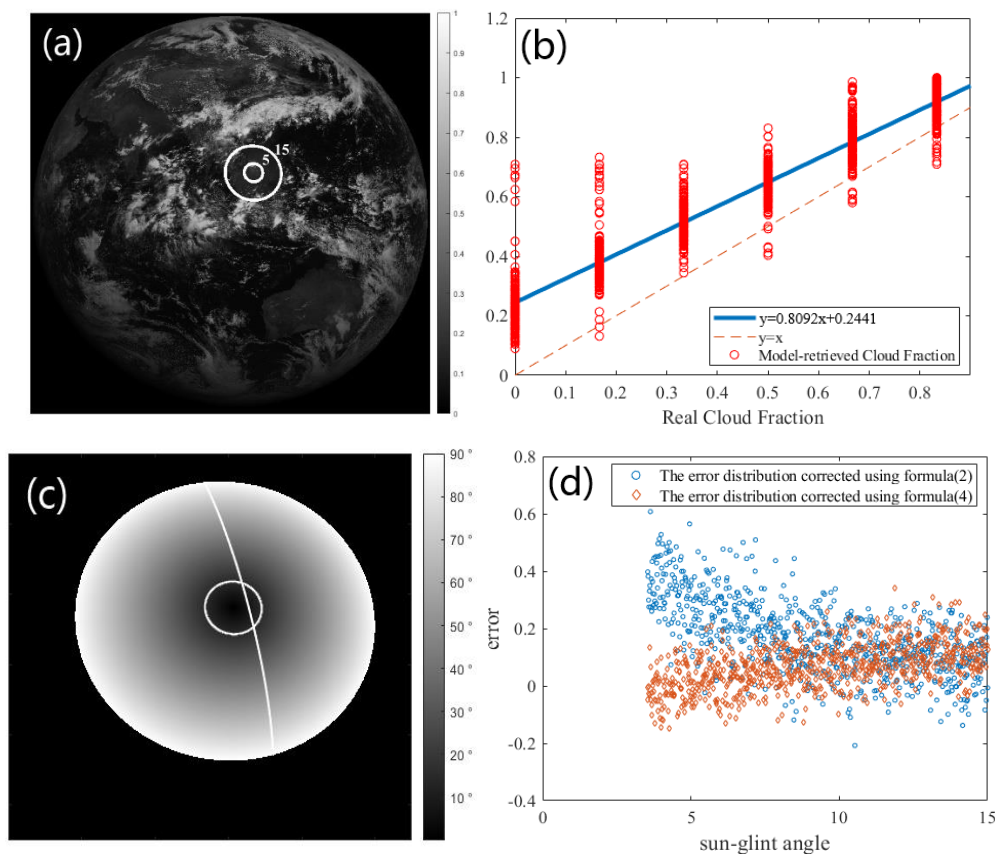
$$414 \quad x = (y - 0.2441)/0.8092 \quad (2)$$

$$415 \quad W_i = \frac{glintangle_i}{\frac{1}{n} \sum_{i=0}^n glintangle_i} \quad (3)$$

$$416 \quad x_i = W_i \left( \frac{y_i - 0.2441}{0.8092} \right) \quad (4)$$

417 Figure 2(d) shows the distribution of errors with respect to SunGlintAngle,  
418 where the blue dots represent the error distribution corrected using formula  
419 (2), and the orange dots represent the error distribution corrected using  
420 formula (4). It can be seen from Figure 2(d) that after correction by formula  
421 (4), the errors in the smaller range of SunGlintAngle are significantly reduced.

422



423

424 **Figure 2:** (a) albedo image of  $0.67\mu\text{m}$  channel (the circles are the contours of the sun-  
425 glint angle), (b) Scatter plot of cloud fraction in sun glint region (The blue line  
426 represents the curve (namely Eq.2) fitted by the least squares method between the  
427 retrievals and the truths.), (c) Distribution of SunGlintAngle and satellite flight track of  
428 CloudSat & Calypso at 4:00 on June 5, 2019, (d) Distribution of cloud fraction retrieval  
429 error with sun-glint angle.

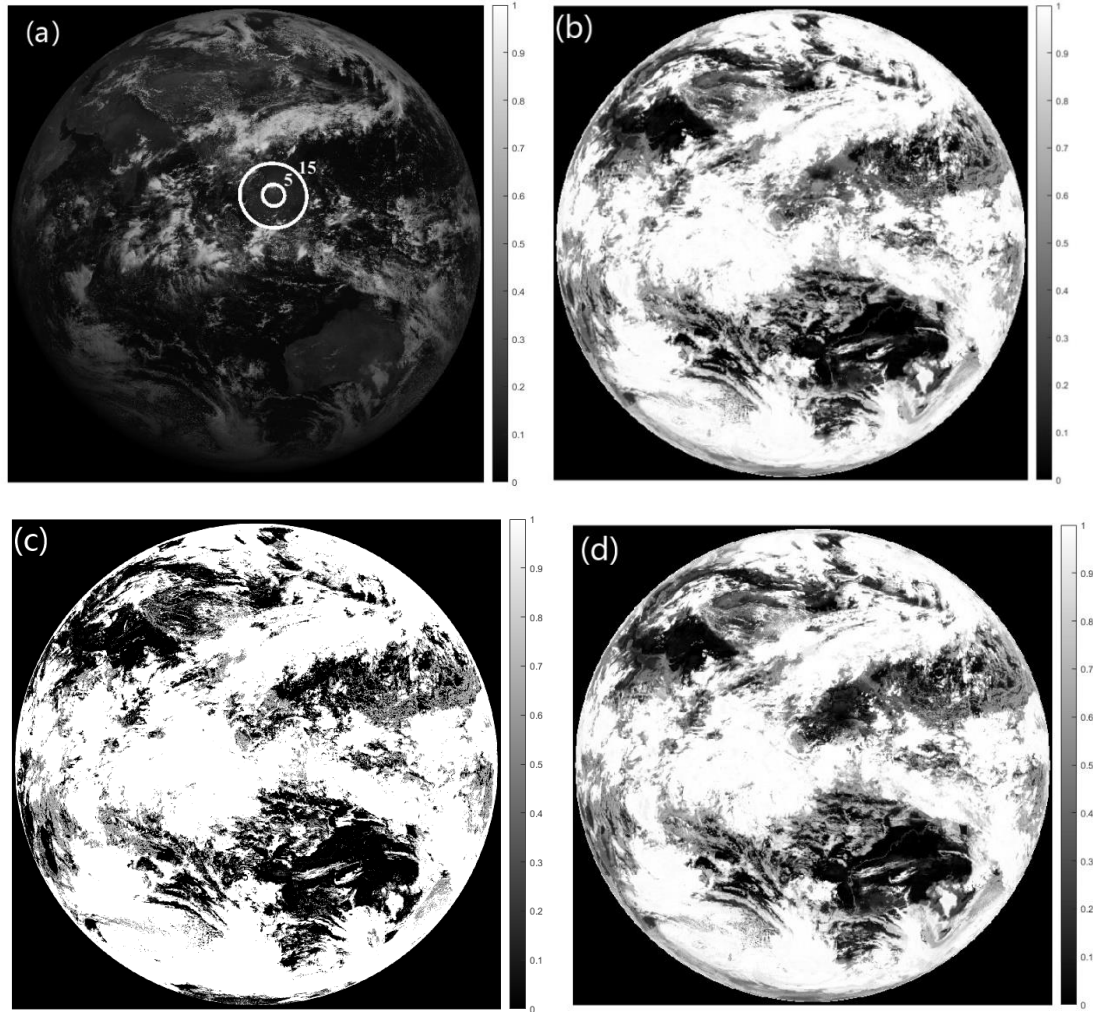
### 430 **3.3 Algorithm universal applicability testing**

431 Although the retrieval model in this article was built based on data from May 2019  
432 due to the limited lifespan of the instrument, how effective is it in real-time FY-4A  
433 AGRI observations and even subsequent FY-4B AGRI applications? The algorithm's  
434 universal applicability was tested using real-time observations from FY-4A and FY-4B  
435 AGRI in 2023.

436 Taking the full-disk observation of FY-4A AGRI at 04:00 (UTC, the same below)  
437 on 1 June 2023 as an example, the radiance observations from 14 channels are initially  
438 fed into the random forest cloud detection model to determine the sky classification

439 (overcast, partly cloudy or clear sky) in each AGRI field. The random forest cloud  
440 fraction retrieval model is utilized to retrieve the cloud fraction in scenes identified as  
441 partly cloudy. Figure 3(a) is the observed albedo at  $0.67\ \mu\text{m}$ , where the circles represent  
442 the contours of the sun glint angle, (b) is the cloud fraction retrievals from random forest  
443 algorithm, (c) is the official operational cloud fraction product and (d) is random forest  
444 cloud fraction retrievals with sun-glint correction. It can be seen from Figure 3 that  
445 many clear-sky scenes are erroneously identified as cloudy by the operational product  
446 and the cloud fraction is generally overestimated with many scenes having a cloud  
447 fraction of 1. The random forest algorithm identifies more regions as clear skies or  
448 partly cloudy than the operational products, matching better with the observations in  
449 the  $0.67\ \mu\text{m}$  albedo image. Brighter regions in the visible image correspond to cloud  
450 cover areas and darker areas represent clear sky conditions. The sun glint region in the  
451 central South China Sea (the circled area in Figure 3(a)) is depicted in Figure 3(b),  
452 where the clear-sky scenes over the ocean are misidentified as partly cloudy by random  
453 forest algorithm due to the increase in observed albedo. Although operational product  
454 in this area also suffers from the impact of unremoved sun glint, it identifies more clear-  
455 sky scenes and the cloud fraction is relatively low. Thus, it is evident that the random  
456 forest algorithm exhibits significant cloud detection and cloud fraction errors in these  
457 sun glint regions. Correction is necessary for the cloud fraction retrievals in the sun  
458 glint region.

459 Figure 3(d) shows the cloud fraction distribution after correction using equation  
460 (9) in the sun glint region., The correction eliminates the influence of sun glint  
461 comparing to the cloud fraction in sun glint area before correction in Figure 3(b). The  
462 scenes misjudged as partly cloudy are corrected to clear sky and match well with the  
463 actual albedo observations in 3(a), which accurately restores the true cloud coverage  
464 over the South China Sea.



465

466

467 **Figure 3:** FY-4A AGRI at 04:00 on 1 June 2023 (a) albedo image of  $0.67\mu\text{m}$  channel  
 468 (the circles are the contours of the sun-glint angle), (b) random forest cloud fraction  
 469 retrieval without sun-glint correction, (c) operational cloud fraction product, (d) random  
 470 forest cloud fraction retrieval with sun-glint correction.

471

472

473

474

475

476

477

Statistical analysis was conducted on the correction effect using samples with sun  
 glint in the training data. The POD and FAR in sun glint area is listed in table 5 and the  
 error is in table 6. It can be seen that after correcting for cloud fraction, the POD for  
 clear skies has increased from 0.0987 to 0.9023. The FAR for partly cloudy has  
 decreased from 0.7943 to 0.0276. Both ME, MAE, and RMSE show significant  
 reductions, and the results after correction outperform operational products.

Table 5 POD and FAR of Cloud Detection in sun glint area

Sky Classification	Operational Product	RF	RF after Correction
-----------------------	------------------------	----	------------------------

	Clear Sky	0.4120	0.0987	0.9023
POD	Partly cloudy	0.7371	0.9663	0.9587
	Overcast	0.8856	0.9845	0.9845
	Clear Sky	0.1229	0.1633	0.0938
FAR	Partly cloudy	0.3332	0.7943	0.0276
	Overcast	0.2983	0.1321	0.1321

478

479

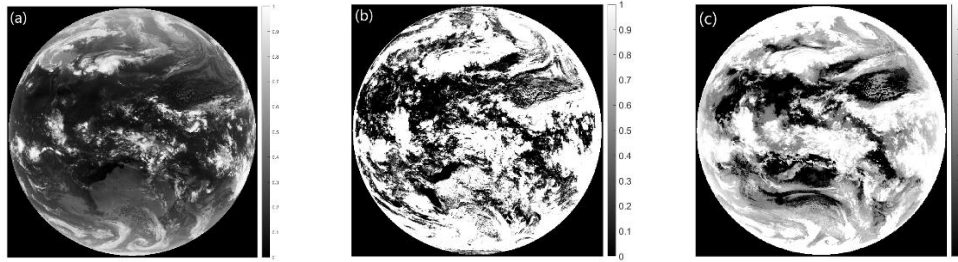
Table 6 cloud fraction Errors in sun glint area

	Operational Product	RF Retrievals	RF after Correction
ME	0.2354	0.1741	0.0670
MAE	0.2511	0.1820	0.0849
RMSE	0.2771	0.2166	0.1041

480

FY-4B launched in 2021 has a total of 15 channels with an additional low-level water vapor channel at 7.42  $\mu\text{m}$  compared to FY-4A. Taking the full-disk observation of FY-4B AGRI at 17:00 on April 18, 2023, as an example, The radiance observation data of the remaining eight channels (near-infrared and infrared channels) except for the 7.42  $\mu\text{m}$  channel and the visible light channels were input into the random forest cloud detection model. Figure 4 (a) shows the brightness temperature distribution observed in the 10.8  $\mu\text{m}$  channel of FY-4B AGRI, (b) represents the operational cloud fraction product for FY-4B AGRI and (c) shows the cloud fraction retrieved by this algorithm. Figure 4 illustrates that the random forest algorithm identifies more regions as clear skies or partly cloudy than the operational products, aligning better with the brightness temperature observations in 10.8  $\mu\text{m}$ . Especially in high latitude regions of the southern hemisphere and areas with strong convection near the equator, the cloud cover provided by operational products is too high and even misjudged. It can be seen that the random forest algorithm is also suitable for cloud fraction retrieval of FY-4B AGRI.

494



495

496 **Figure 4:** FY-4B AGRI at 17:00 on 18 April 2023, (a) brightness temperature of  $10.8\mu\text{m}$   
 497 channel, (b) operational cloud fraction product, (c) random forest cloud fraction  
 498 retrieval.

#### 499 **4 Conclusion**

500 This paper used the random forest and multi-layer perceptron (MLP) algorithms  
 501 to retrieve cloud fraction from FY-4A AGRI full-disk Level-1 radiance observation data,  
 502 and verified the accuracy of the algorithms using the Cloudsat & Calypso active remote  
 503 sensing satellite's 2B CLDCLASS-LIDAR cloud fraction product. The following  
 504 conclusions were drawn:

505 (1) The random forest and MLP algorithms performed well in cloud detection and  
 506 cloud fraction retrieval tasks, and their accuracy was higher than that of operational  
 507 products. The accuracy of cloud detection can reach over 93%, and the error of cloud  
 508 fraction retrieval is close to zero. Compared with the MLP algorithm, the RF algorithm  
 509 has a slightly higher accuracy in cloud detection, and a slightly lower error in cloud  
 510 fraction retrieval, showing better performance.

511 (2) At night, the classification accuracy is lower than during the day due to the lack  
 512 of observations in the visible channel of AGRI, resulting in higher cloud fraction errors  
 513 at night.

514 (3) The accuracy of identifying partly cloudy scenes is lower than that of  
 515 identifying clear sky and overcast scenes for both RF and MLP algorithms. Scenes with  
 516 very low cloud fraction (0.16) are often misclassified as clear sky, while scenes with  
 517 high cloud fraction (0.83) are often misclassified as overcast.

518 (4) The sun-glint area cloud fraction correction curve, fitted with SunGlintAngle  
519 as the weight, greatly improves the accuracy of cloud fraction retrieval and reduces the  
520 misclassification rate of clear sky scenes as partly cloudy or partly cloudy scenes as  
521 overcast due to increased reflectance.

522

### 523 *Data availability*

524 FY-4A AGRI data is available at <http://satellite.nsmc.org.cn> and the 2B-CLDCLASS-  
525 LIDAR data at <https://www.icare.univ-lille.fr/data-access/data-archive-access/>

526

### 527 *Author contributions*

528 JX: Formal analysis, Methodology, Software, Visualization and Writing – original draft  
529 preparation. LG: Conceptualization, Data curation, Funding acquisition, Supervision,  
530 Validation and Writing – review & editing.

531

### 532 *Competing interests*

533 The contact author has declared that none of the authors has any competing interests.

534

### 535 *Disclaimer*

### 536 *Acknowledgements*

537 Funding: This work was supported by the National Natural Science Foundation of  
538 China under grant no. 41975028.

539 *We acknowledge the High Performance Computing Center of Nanjing University of*  
540 *Information Science & Technology for their support of this work.*

### 541 **References**

542 Amato, U., Antoniadis, A., Cuomo, V., Cuttillo, L., Franzese, M., Murino, L., Serio,  
543 C.: Statistical cloud detection from SEVIRI multispectral images, *Remote Sensing of*  
544 *Environment*, 112, 750–766, <https://doi.org/10.1016/j.rse.2007.06.004>, 2008.

545 Baum, B., Trepte Q.: A Grouped Threshold Approach for Scene Identification in  
546 AVHRR Imagery, *Journal of Atmospheric & Oceanic Technology*, 16, 793-800,  
547 [https://doi.org/10.1175/1520-0426\(1999\)016<0793:AGTAFS>2.0.CO;2](https://doi.org/10.1175/1520-0426(1999)016<0793:AGTAFS>2.0.CO;2), 1999.

548 Breiman L.1999. Random Forests-Random Features [J]. *Machine Learning*.45(1): 5-32.  
549 Breiman, L. Random Forests. *Machine Learning* **45**, 5–32 (2001).  
550 [doi.org/10.1023/A:1010933404324](https://doi.org/10.1023/A:1010933404324)

551 Chai D ,Huang J ,Wu M , et al.Remote sensing image cloud detection using a shallow  
552 convolutional neural network[J].*ISPRS Journal of Photogrammetry and Remote*  
553 *Sensing*,2024,20966-84.

554 Merchant, C.J., Harris, A.R., Maturi, E., Maccallum S.: Probabilistic physically based  
555 cloud screening of satellite infrared imagery for operational sea surface temperature  
556 retrieval, *Quarterly Journal of the Royal Meteorological Society*, 131, 2735-2755,  
557 <https://doi.org/10.1256/qj.05.15>, 2005.

558 Gao, J., Jing, Y.: Satellite Remote Sensing Cloud Detection Method Based on Fully  
559 Convolutional Neural Network,*Infrared Technology*, 41, 607-615, 2019.

560 Gomez-Chova, L., Camps-Valls, G., Amoros-Lopez, J., Guanter, L., Alonso, L., Calpe,  
561 J., Moreno, J.: New Cloud Detection Algorithm for Multispectral and Hyperspectral  
562 Images: Application to ENVISAT/MERIS and PROBA/CHRIS Sensors, *IEEE*  
563 *International Symposium on Geoscience and Remote Sensing*, 2757–  
564 2760, doi:10.1109/igarss.2006.709, 2006.

565 Kay, S., Hedley, J., Lavender, S.: Sun Glint Correction of High and Low Spatial  
566 Resolution Images of Aquatic Scenes: a Review of Methods for Visible and Near-  
567 Infrared Wavelengths, *Remote Sensing*, 1, 697-730,  
568 <https://doi.org/10.3390/rs1040697>, 2009.

569 Kegelmeyer, W.P.J.: Extraction of cloud statistics from whole sky imaging  
570 cameras,1994.

571 Kong, Y.-L., Huang, Q., Wang, C., Chen, J., Chen, J., & He, D. (2018). Long Short-  
572 Term Memory Neural Networks for Online Disturbance Detection in Satellite Image  
573 Time Series. *Remote Sensing*, 10(3), 452. doi:10.3390/rs10030452

574 Mace, G. G., R. Marchand, Q. Zhang, et al. (2007). CloudSat Project: Level 2 Radar-  
575 Lidar GEOPROF product process description and interface control document. Jet  
576 Propulsion Laboratory.

577 Pan, C., Xia B., Chen, Y.: Research on MODIS Cloud Detection Algorithms Based on  
578 Fuzzy Clustering, *Microcomputer Information*, 25, 124-125+131, 2009.

579 Yan J, Guo X, Qu J, Han M. An FY-4A/AGRI cloud detection model based on the naive  
580 Bayes algorithm. *Remote Sensing for Natural Resources*, 34(3): 33-42. doi:  
581 10.6046/zrzyyg.2021259. 2022

582 Rossow, W. B., Leonid, C.G.: Cloud detection using satellite measurements of infrared  
583 and visible radiances for ISCCP. *Journal of Climate*, 12, 2341-2369,  
584 [https://doi.org/10.1175/1520-0442\(1993\)006<2341:CDUSMO>2.0.CO;2](https://doi.org/10.1175/1520-0442(1993)006<2341:CDUSMO>2.0.CO;2), 1993.

585 R.O. Duda, P.E. Hart, and D.G. Stork, *Pattern Classification*, New York: John Wiley &



586 Sons, 2001, pp. xx + 654, ISBN: 0-471-05669-3. *Journal of Classification* 24, 305–  
587 307 (2007). <https://doi.org/10.1007/s00357-007-0015-9>  
588 Solvsteen, C.: Correlation based cloud-detection and an examination of the split-  
589 window method, *Proceedings of SPIE - The International Society for Optical*  
590 *Engineering*, 86-97, 1995.  
591 Wang, Z.: CloudSat Project: CloudSat 2B-CLDCLASS-LIDAR product process  
592 description and interface control document, *Jet Propulsion Laboratory*, 2019.  
593 Yan, J., Guo, X., Qu, J.: An FY-4A/AGRI cloud detection model based on the naive  
594 Bayes algorithm, *Remote Sensing for Natural Resources*, 34, 33-42, 2022.  
595 Zhang, W., He, M., Mak, M.W.: Cloud detection using probabilistic neural networks,  
596 *Geoscience and Remote Sensing Symposium*, IEEE 2373-2375, 2001.  
597 Zhang, Y., William, B. R., Andrew, A. L., Valdar, O., Michael, I. M.: Calculation of  
598 radiative fluxes from the surface to the top of atmo- sphere based on ISCCP and other  
599 global data sets: Refine- ments of the radiative transfer model and the input data,  
600 *Journal of Geophysical Research Atmospheres*, 109, 1-27,  
601 <https://doi.org/10.1029/2003JD004457>, 2004.