

20 identified as partly cloudy. The 2B-CLDCLASS-LIDAR cloud fraction product from
21 Cloudsat& CALIPSO active remote sensing satellite is employed as the truth to assess
22 the accuracy of the retrieval algorithm. Comparison with the operational AGRI level 2
23 cloud fraction product is also conducted at the same time. The results indicate that both
24 the Random Forest (RF) and Multi-Layer Perceptron (MLP) cloud detection models
25 achieved high accuracy, surpassing that of operational products. However, both
26 algorithms demonstrated weaker discrimination capabilities for partly cloudy
27 conditions compared to clear sky and overcast situations. Specifically, they tended to
28 misclassify fields of view with low cloud fractions (e.g., cloud fraction = 0.16) as clear
29 sky and those with higher cloud fractions (e.g., cloud fraction = 0.83) as overcast.
30 Between the two models, RF exhibited higher overall accuracy. Both RF and MLP
31 models performed well in cloud fraction retrieval, showing lower mean error (ME),
32 mean absolute error (MAE), and root mean square error (RMSE) compared to
33 operational products. The ME for both RF and MLP cloud fraction retrieval models was
34 close to zero, while RF had slightly lower MAE and RMSE than MLP. During daytime,
35 the high reflectance in sun-glint areas led to larger retrieval errors for both RF and MLP
36 algorithms. However, after correction, the retrieval accuracy in these regions improved
37 significantly. At night, the absence of visible light observations from the AGRI
38 instrument resulted in lower classification accuracy compared to daytime, leading to
39 higher cloud fraction retrieval errors during nighttime.

40 **Key words:** Cloud detection; cloud fraction retrieval; FY-4A AGRI; CloudSat &

41 CALIPSO; machine learning; deep learning.

42 **Introduction**

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

54 Fengyun-4A (FY-4A) is a comprehensive atmospheric observation satellite
55 launched by China in 2016. The uploaded AGRI (Advanced Geosynchronous Radiation
56 Imager) has 14 channels and captures full-disk observation every 15 minutes. In
57 addition to observing clouds, water vapor, vegetation and the Earth's surface, it also
58 possesses the capability to capture aerosols and snow. Moreover, it can clearly
59 distinguish different phases and particle size of clouds and obtain high- to mid-level

60 water vapor content. It is particularly suitable for cloud detection due to its
61 simultaneous use of visible, near-infrared, and long-wave infrared channels for
62 observation with 4km spatial resolution.

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

80 The other category of cloud detection algorithms is based on statistical probability

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

90 More and more machine learning algorithms are being utilized by researchers in
91 cloud detection studies with the development of machine learning. For instance, the
92 probabilistic neural networks, especially radial basis function networks was used for
93 AVHRR cloud detection (Zhang, et al., 2001). The utilization of convolutional neural
94 network methods (Hu, et al., 2020) offers important perspectives for cloud detection
95 research.

96 Currently, there is limited research literature on cloud detection and cloud fraction
97 retrieval algorithms for FY-4A/4B AGRI. The operational cloud fraction product of FY-
98 4A AGRI utilized a threshold method with 4 km spatial resolution. Differences in
99 climatic and environmental factors lead to varying albedo and brightness temperature
100 observations for the instrument at different times and locations. Therefore, the choice
101 of thresholds is easily influenced by factors such as season, latitude and land surface

102 type (Gao and Jing, 2019). Using multiple sets of thresholds for discrimination would
103 significantly slow down the cloud detection process. Moreover, most algorithms focus
104 solely on cloud detection, which classified the observed scenes into cloud or clear-sky
105 without providing the specific cloud fraction information for the scenes. The use of
106 active remote sensing instruments carried by Cloudsat & Calypso is not influenced by
107 thresholds when retrieving cloud fraction, enabling a more accurate cloud fraction
108 retrieval. However, due to Cloudsat & Calypso being polar-orbiting satellites, the cloud
109 fraction over the full disk cannot be obtained. Utilizing the Cloudsat & Calypso Level
110 2 product 2B-CLDCLASS-LIDAR as the reference truth, a random forest model trained
111 based on FY4A AGRI full disk radiation data can address the shortcomings of threshold
112 methods and achieve a high accuracy of cloud fraction over the full disk.

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

119 **1 Research Data and Preprocessing**

120 *1.1 FY-4A data*

121 FY-4A was successfully launched on December 11, 2016. Starting from May 25, 2017,
122 FY-4A drifted to a position near the main business location of the Fengyun
123 geostationary satellite at 104.7 degrees east longitude on the equator. Its successful
124 launch marked the beginning of a new era for China's next-generation geostationary
125 meteorological satellites as an advanced comprehensive atmospheric observation
126 satellite. The Advanced Geosynchronous Radiation Imager (AGRI), one of the main
127 payloads of the Fengyun-4 series geostationary meteorological satellites, can perform
128 large-disk scans and rapid regional scans at a minute level. It has 14 observation
129 channels in total with the main task of acquiring cloud images. The channel parameters
130 and main uses of AGRI are detailed in Table 1
131 (<https://www.nsmc.org.cn/nsmc/cn/instrument/AGRI.html>). The first six visible light
132 channels have no values at night, meaning that channels with a central wavelength less
133 than or equal to $2.225\mu\text{m}$ are unavailable during nighttime. FY-4A AGRI data was
134 downloaded from the official website of the China national satellite meteorological
135 center (<http://satellite.nsmc.org.cn>), including level-1 full disk radiation observation
136 data preprocessed through quality control, geolocation and radiation calibration as well
137 as level-2 cloud fraction product (CFR). The spatial resolution of these data is all 4 km
138 at nadir and the temporal resolution is 15 minutes.

Table 1 FY-4A AGRI channel parameters

Channel Number	Band Range / μm	Central Wavelength / μm	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
13	11.50 ~ 12.50	12.0	4	surface and cloud-top temperatures
14	13.2 ~ 13.8	13.5	4	cloud-top height

140

141 *1.2 CloudSat & Calipso Cloud Product*

142 CALIPSO (Cloud-Aerosol Lidar and Infrared Pathfinder Satellite Observations)

143 is a satellite jointly launched by NASA and CNES (the French National Center for

144 Space Studies) in 2006. It is a member of the A-Train satellite observation system.

145 CALIPSO is equipped with three payloads, among which CALIOP (the Cloud and

146 Aerosol Lidar with Orthogonal Polarization) is a primary observational instrument.

147 Observing with dual wavelengths (532 nm and 1064 nm) CALIOP can provide high-

148 resolution vertical profiles of clouds and aerosols with 30 m vertical resolution. As the

149 first satellite designed to observe global cloud characteristics in a sun-synchronous orbit

150 CloudSat is also among NASA's A-Train series satellites. The CPR (Cloud Profile
151 Radar) installed on it operates at 94 GHz millimeter-wave and is capable of detecting
152 the vertical structure of clouds and providing vertical profiles of cloud parameters. The
153 scanning wavelengths of CPR and CALIOP are different. CALIOP is capable of
154 observing the top of mid-to-high level clouds, whereas CPR can penetrate optically
155 thick clouds. Combining the strengths of these two instruments enables the acquisition
156 of precise and detailed information on cloud layers and cloud fraction.

157 The joint level 2 product 2B-CLDCLASS-LIDAR is mainly utilizing in this study.
158 It provides the cloud fraction at different heights with horizontal resolution 2.5 km
159 (along-track) \times 1.4 km (cross-track) through combining the observations from CPR and
160 CALIOP. Since the two instruments have different spatial domain such as vertical
161 resolution, spatial resolution and spatial frequency, the spatial domain of the output
162 products is defined in terms of the spatial grid of the CPR. In the algorithm, the cloud
163 fraction is calculated using a weighted scheme based on the spatial probability of
164 overlap between the radar and lidar observations. The calculation of the lidar cloud
165 fraction within a radar footprint is represented by the equation 1 (Mace, G. G., et al,
166 2007):

$$167 \quad C_l = \frac{\sum_{i=1}^{\# \text{ of lidar obs}} w_i \delta_i}{\sum_{i=1}^{\# \text{ of lidar obs}} w_i} \quad (1)$$

168 Where:

169 C_l represents the lidar cloud fraction within a radar footprint.

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

171 δ_i indicates the lidar hydrometeor occurrence, where a value of 1 signifies the
172 presence of hydrometeor and 0 indicates the absence.

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

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

180 ***1.3 Establishment of Training Data***

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

187 The FY-4A AGRI and 2B-CLDCLASS-LIDAR data with a spatial difference
188 between fields of view within 1.5 km and a time difference within 15 minutes are
189 spatiotemporal matched. To make the 2B-CLDCLASS-LIDAR cloud fraction data
190 collocated within AGRI pixels more effective, at least two 2B-CLDCLASS-LIDAR

191 pixels are required within each AGRI field of view. The cloud fraction average of these
192 pixels is used as the cloud fraction for that AGRI pixel. However, the errors in the
193 matched dataset are unavoidable. The AGRI scanning method operates from left to right
194 and top to bottom. Each complete scan of the full disk takes 15 minutes and generates
195 a dataset. It is impossible to determine the exact moment of a specific point within the
196 full disk. This limits the time range for matching datasets to within 15 minutes.
197 However, in areas with higher wind speeds, clouds can move a significant distance
198 within that 15-minute window. Therefore, errors arising from timing issues cannot be
199 avoided.

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

207 The algorithm includes two steps: the cloud detection is conducted firstly for each
208 AGRI field of view to identify whether it is clear sky, partly cloudy or overcast within
209 the observation field. Then the cloud fraction is retrieved for the scene identified as
210 partly cloudy. So the training data include A dataset used for cloud detection and B
211 dataset for cloud fraction retrieval. The input variables in A dataset are the FY-4A

212 AGRI level-1 radiative observations from 14 channels and the output variable is the
213 temporally and spatially matched 2B-CLDCLASS-LIDAR cloud detection label. The
214 output is categorized into three types: overcast, partly cloudy and clear sky with values
215 1, 2 and 3 respectively. The cloud fraction product from 2B-CLDCLASS-LIDAR
216 consists of discrete values: 0, 0.16, 0.33, 0.50, 0.66, 0.83, and 1. Here, 0 indicates clear
217 sky, values from 0 to 1 represent varying cloud fractions for partly cloudy conditions,
218 and 1 signifies overcast. To ensure the balance and representativeness of the samples,
219 the proportions of different cloud fraction samples in dataset A are set at 5:1:1:1:1:5.
220 Regarding the samples for partly cloudy type in dataset A, the collocated 2B-
221 CLDCLASS-LIDAR cloud fraction products serve as output labels for cloud fraction
222 retrieval model B. The input of training dataset B remains the FY-4A AGRI level-1
223 radiative observations.

224 Due to the instrument's limited lifespan, only 2B-CLDCLASS-LIDAR data up to
225 August 2019 can be obtained. The sample time range used in this paper is from August
226 2018 to July 2019. Five days were randomly selected each month as daytime samples
227 and five days as nighttime samples. A total of 120 days of time and space matched FY-
228 4A AGRI full-disk observations and 2B-CLDCLASS-LIDAR data were used as
229 training and testing samples. Among them, 80% of the data was used for training, and
230 20% was used for testing. The total number of daytime samples in dataset A is 91,073,
231 while dataset B contains 30,358 samples. The total number of nighttime samples in
232 dataset A is 95,493, and dataset B includes 31,831 samples.

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

236

237 **2 Algorithms**

238 Our preliminary experiments involved multiple algorithms, including LibSvm,
239 MLP, BP neural network, and Random Forest. These experiments highlighted that,
240 among the baselines, Random Forest and MLP achieved the highest overall accuracy.
241 For this reason, we selected them to perform additional experiments.

242 **2.1 Random Forest (RF)**

243 This algorithm integrates multiple trees based on the Bagging idea of ensemble
244 learning, with the basic element being the decision tree ([Breiman, 1999](#)). When building
245 a decision tree, N sets of independent and dependent variables are randomly sampled
246 with replacement from the original training samples to create a new training sample set;
247 m variables are randomly sampled without replacement from all independent variables,
248 the dependent variable data is split into two parts using the selected variables, and the
249 purity of the subsets is calculated for each split method. The variable utilized by the
250 split method with the highest purity is used to partition the data, completing the decision
251 at that node. This process of binary splitting continues to grow the decision tree until

252 stopping criteria are met, completing the construction of a single decision tree. These
253 steps are repeated N_{tree} times to build a random forest model consisting of N_{tree}
254 decision trees (Breiman, 2001). Random Forest adopts ensemble algorithms, with the
255 advantage of high accuracy. It can handle both discrete and continuous data, without
256 the need for normalization, making it more efficient compared to other algorithms.

257 ***2.2 Multilayer Perceptron (MLP)***

258 This algorithm consists of a fully connected artificial neural network (Duda, et al.,
259 2001). The classifier/regressor takes feature vectors or tensors as input. The input is
260 mapped through multiple fully connected hidden layers containing hidden weights,
261 which produce classifications/regressions at the output layer. A nonlinear activation
262 function (such as sigmoid or rectified linear unit (ReLU)) is applied in each hidden
263 layer to facilitate a nonlinear model. For classifiers, the output of the final hidden layer
264 is combined and passed through a softmax function to generate class predictions. The
265 model's weights are trained in a supervised manner, utilizing stochastic gradient descent
266 and backpropagation to achieve the desired classification/regression.

267 ***2.3 Hyperparameters***

268 In this paper, a total of eight models were established, including daytime/nighttime
269 random forest classification/regression models and daytime/nighttime MLP
270 classification/regression models. For the random forest, we first conducted experiments

271 using the following Hyperparameters ranges: Trees: [200, 300, 400, 500, 600,700],
272 minleaf: [1, 2, 5, 10], criterion: [Gini, entropy]. Ultimately, the best selections were: (1)
273 Daytime RF classification model: Trees=500, minleaf=1, criterion=gini; (2) Nighttime
274 RF classification model: Trees=600, minleaf=1, criterion=gini; (3) Daytime RF
275 regression model: Trees=400, minleaf=1, criterion=gini; (4) Nighttime RF regression
276 model: Trees=500, minleaf=1, criterion=gini.

277 For the MLP, experiments were conducted using the following hyperparameter
278 ranges: Hidden layer size: [2,3,4,5,6,7,8,9], Hidden layer neuron count:
279 [8,16,32,64,128], Activation hyperparameter: [logistic, tanh, relu], MaxEpochs:
280 [30,50,100], MiniBatchSize: [300,400,....,1500,1600], Solver hyperparameter: [lbfgs,
281 sgd, adam]. The optimal parameters found are as follows: (1) MLP classification model
282 for daytime: hidden layer size = 5, MiniBatchSize = 1500. (2) MLP classification model
283 for nighttime: hidden layer size = 7, MiniBatchSize = 800. (3) MLP regression model
284 for daytime: hidden layer size = 4, MiniBatchSize = 600. (4) MLP regression model for
285 nighttime: hidden layer size = 6, MiniBatchSize = 500. All four models have hidden
286 layer neuron count = 64, activation = relu, MaxEpochs = 50, solver = adam,
287 InitialLearnRate = 0.01, LearnRateSchedule = piecewise, LearnRateDropFactor = 0.1,
288 LearnRateDropPeriod = 10.

289 **3 Results and Analysis**

290 To assess the accuracy and stability of the retrieval model, two types of validation
291 methods are utilized. One way involves a direct comparison from images, qualitatively
292 comparing the model's retrieval results and official cloud fraction products with AGRI
293 observed cloud images. Another approach uses 2B-CLDCLASS-LIDAR as the ground
294 truth and introduces five parameters for quantitative comparison: recall, false alarm rate
295 (FAR), mean error (ME), mean absolute error (MAE), and root mean square error
296 (RMSE). To evaluate the ability of operational products, RF, and MLP cloud detection
297 models to distinguish overcast, partly cloudy, and clear sky, the recall is calculated using
298 the formula $POD=TP/(TP+FN)$, and the false alarm rate is calculated using the formula
299 $FAR=FP/(TP+FP)$. Taking the overcast scene as an example, TP represents the number
300 of correctly identified overcast conditions, FN represents the number of overcast
301 conditions misidentified as partly cloudy or clear sky, and FP represents the number of
302 clear sky or partly cloudy conditions misidentified as overcast. When assessing the
303 accuracy of operational products and cloud fraction models for the cloud fraction
304 retrieval results of partly cloudy scenes, mean error (ME), mean absolute error (MAE),
305 and root mean square error (RMSE) are used.

306 ***3.1 Objective Analysis of Cloud Fraction Retrievals***

307 First, using the 2B-CLDCLASS-LIDAR cloud fraction product as the ground truth,

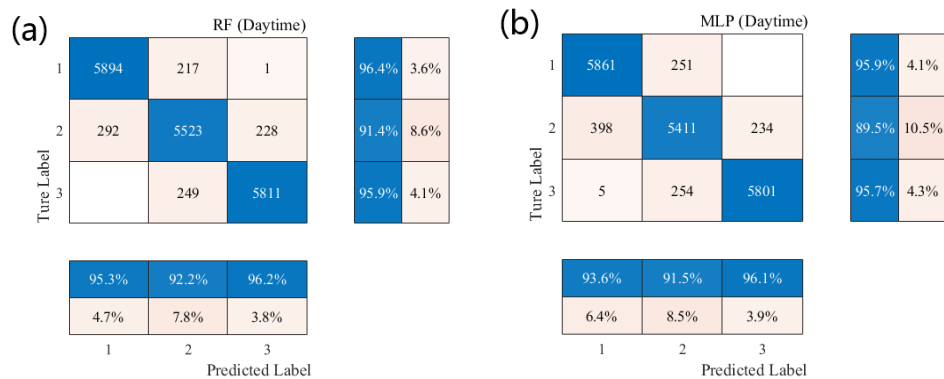
308 we calculated the accuracy of the operational cloud detection products. The results are
 309 shown in Table 2. The samples used for this statistic are the same as those for testing
 310 the model below (20% of dataset A).

311 Table 2: Recall Rate and FAR of Operational Cloud Detection Products

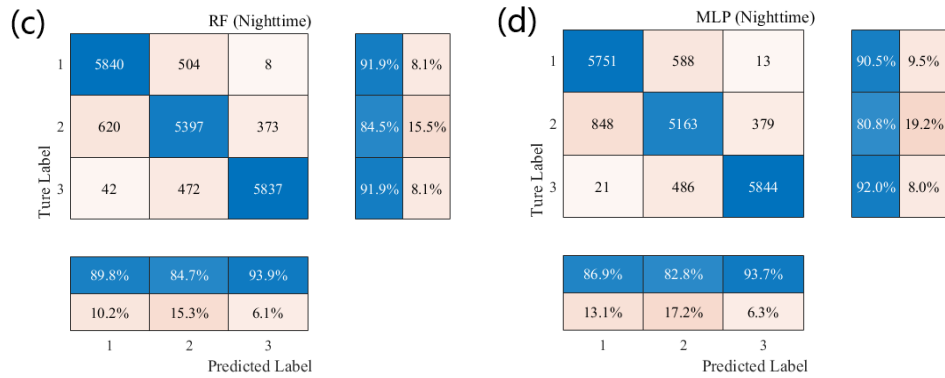
	Sky Classification	Daytime Product	Nighttime Product
	Clear Sky	0.6359	0.5781
POD	Partly cloudy	0.7174	0.7449
	Overcast	0.7736	0.7384
	Clear Sky	0.1778	0.0934
FAR	Partly cloudy	0.1819	0.2117
	Overcast	0.2499	0.2683

312 Based on the cloud detection model trained above, cloud detection experiments
 313 were conducted using the test samples from Dataset A. The time-space matched 2B
 314 CLDCLASS-LIDAR cloud fraction product served as the ground truth to assess the
 315 accuracy of cloud detection. Figure 1 shows the results: (a) Random Forest model
 316 results during the day, (b) MLP model results during the day, (c) Random Forest model
 317 results during the night, and (d) MLP model results during the night. The x-axis
 318 represents the model predictions, while the y-axis represents the ground truth. A value
 319 of 1 on both axes indicates clear skies, 2 indicates partly cloudy, and 3 indicates overcast.
 320 The blue area on the right side of each plot shows the recall rate for each type, while
 321 the light-colored area at the bottom represents the False Alarm Rate (FAR). During the

322 day, the Random Forest model achieved an overall accuracy of 94.2%, while the MLP
 323 model had an overall accuracy of 93.4%. The Random Forest model exhibited slightly
 324 higher recall rates for clear skies, partly cloudy, and overcast conditions compared to
 325 the MLP model, and its FAR was lower as well. Both models performed poorly in
 326 recognizing partly cloudy conditions, as the models tended to classify true cloud
 327 fractions of 0.16 as clear skies and those of 0.83 as overcast. At night, the Random
 328 Forest model achieved an overall accuracy of 89.4%, while the MLP model had an
 329 accuracy of 87.7%. The Random Forest model had higher recall rates for clear skies
 330 and partly cloudy conditions compared to the MLP, while the recall rates for overcast
 331 conditions were similar for both models. The FAR for the Random Forest model was
 332 lower than that of the MLP. Overall, both the Random Forest and MLP models showed
 333 higher classification accuracy for clear skies, partly cloudy, and overcast conditions
 334 compared to operational products, with the Random Forest model performing better.



335



336

337

338

339

Figure 1 Model Cloud Detection Accuracy: (a) Daytime RF, (b) Daytime MLP, (c) Nighttime RF, (d) Nighttime MLP (In the axis, 1 represents clear sky, 2 represents partly cloudy, and 3 represents overcast.)

341

342

343

344

345

346

347

348

349

350

351

352

353

Based on the previous model's assessment of the field of view as partly cloudy, the cloud fraction in this AGRI field of view is retrieved using the cloud fraction model established earlier. For model evaluation, both the operational product and the 2B-CLDCLASS-LIDAR cloud fraction product are classified as partly cloudy, with the matched 2B-CLDCLASS-LIDAR cloud fraction product considered as the ground truth. The average error, mean absolute error, and root mean square error for both daytime and nighttime operational products (Table 3) and cloud fraction model retrieval (Table 4) are calculated. It can be observed that the average errors of both models are close to 0 during both daytime and nighttime. The errors are smaller during the day than at night, with the RF model exhibiting lower errors than the MLP model. In summary, the errors of both models are smaller than those of the operational products, and the RF model performs better in the cloud fraction retrieval task.

Table 3: Errors of Operational Product Cloud Fraction

	Daytime Operational	Nighttime Operational
	Product	Product
ME	0.1987	0.2121
MAE	0.2279	0.2441
RMSE	0.2776	0.2938

354

Table 4: Model Retrieval Error

	Daytime	Daytime	Nighttime	Nighttime
	RF	MLP	RF	MLP
ME	0.0006	-0.0009	-0.0028	-0.0032
MAE	0.1011	0.1053	0.1221	0.1322
RMSE	0.1285	0.1332	0.1510	0.1623

355

Based on the experiments mentioned above, the performance of RF in cloud

356

detection and cloud fraction retrieval slightly outperforms that of MLP. Therefore,

357

subsequent experiments will utilize the RF algorithm.

358

3.2 Cloud fraction correction in sun glint regions

359

Sun glint refers to the bright areas created by the reflection of sunlight to the

360

sensors of observation systems (satellites or aircrafts). This phenomenon usually occurs

361

on extensive water surfaces, such as oceans lakes or rivers. This specular reflection of

362

sunlight will cause an increase in the reflected solar radiation received by onboard

363

sensors, manifested as an enhancement of white brightness in visible images. The

364 increase in visible channel observation albedo will affect various subsequent
365 applications of data, including cloud detection and cloud cover retrieval, etc.

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

376 To correct the cloud fraction in the sun-glint areas, we first identified the fields of
377 view (FOVs) where sun-glint occurred during FY-4A AGRI observations from August
378 2018 to July 2019, totaling 1,476 FOVs. Subsequently, a direct least squares fitting
379 was conducted between the retrieved cloud fraction and the collocated 2B-
380 CLDCLASS-LIDAR cloud fraction (ground truth). The scatter plot is illustrated in
381 Figure 2(b), where x-axis is the 2B-CLDCLASS-LIDAR cloud fraction and y-axis is
382 the model-retrieved cloud fraction. The blue line represents the curve (namely Eq.2)
383 fitted by the least squares method between the retrievals and the truths. The thin dash
384 line is the $x=y$ line. It is evident that the retrieved cloud fraction is generally slightly

385 overestimated.

386 Taking observations at 04:00 on 5 June 2019 as an example, Figure 2(c) presents
387 the distribution of SunGlintAngle and the flight trajectory of the Cloudsat&Calypso
388 satellite. White circles denote the sun glint region with SunGlintAngle $< 15^\circ$ and the
389 white line represents the satellite flight track. As depicted in the figure, the majority of
390 Cloudsat&Calypso flight trajectories do not pass through the central position of sun
391 glint area but instead traverse locations with larger SunGlintAngle values. The
392 intensity of sun glint effect decreases with the increase of SunGlintAngle. This
393 suggests that the true values for spatial and temporal matching mostly do not fall within
394 the strongest sun glint region. From Figure 2(d), it can be seen that the impact of sun
395 glint becomes stronger as SunGlintAngle decreasing, which results in a higher
396 observation albedo. This further leads to the overestimated cloud fraction values in the
397 retrieval. It is evident that the cloud fraction error is related to the value of
398 SunGlintAngle and this influence is not considered in Eq. (2). Directly applying
399 equation (2) to correct the cloud fraction retrievals would result in a too small correction
400 intensity for the FOVs near the center of sun glint and an excessively large correction
401 intensity for the FOVs in the Sun-glint edge region (even erroneous clear sky may
402 appear). Considering this, a correction formula (3)-(4) using SunGlintAngle as weight
403 is introduced, where W_i represents the angle weight for a certain pixel i in the sun glint
404 region, n is the number of pixels within the SunGlintAngle $< 15^\circ$ range, y_i is the initial
405 model retrieval of cloud cover for the field of view i and x_i is the final corrected cloud

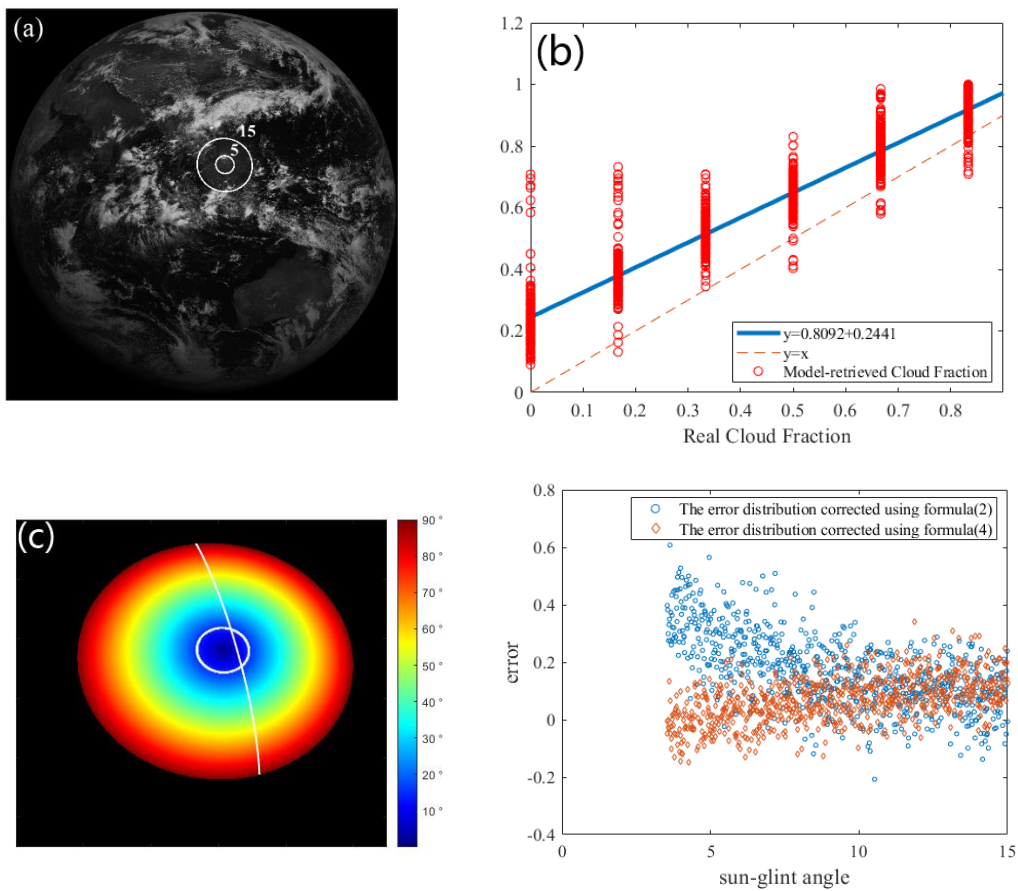
406 fraction.

$$407 \quad x = (y - 0.2441)/0.8092 \quad (2)$$

$$408 \quad W_i = \frac{glintangle_i}{\frac{1}{n} \sum_{i=0}^n glintangle_i} \quad (3)$$

$$409 \quad x_i = W_i \left(\frac{y_i - 0.2441}{0.8092} \right) \quad (4)$$

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



415

416

417 **Figure 2** (a) albedo image of 0.67 μ m channel (the circles are the contours of the sun-

418 glint angle), (b) Scatter plot of cloud fraction in sun glint region (The blue line
419 represents the curve (namely Eq.2) fitted by the least squares method between the
420 retrievals and the truths.), (c) Distribution of SunGlintAngle and satellite flight track of
421 CloudSat & Calypso at 4:00 on June 5, 2019, (d) Distribution of cloud fraction retrieval
422 error with sun-glint angle.

423 *3.3 Algorithm universal applicability testing*

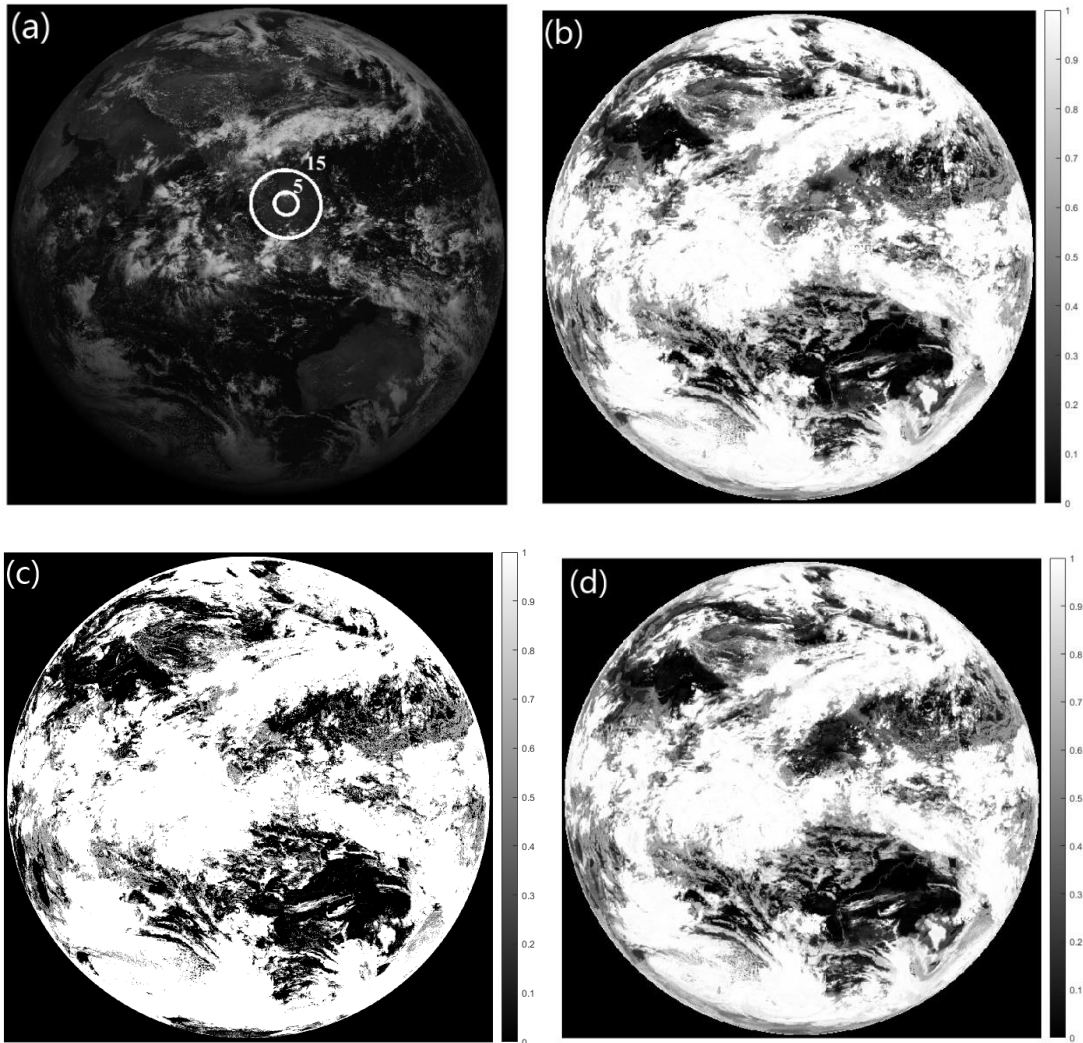
424 Although the retrieval model in this article was built based on data from May 2019
425 due to the limited lifespan of the instrument, how effective is it in real-time FY-4A
426 AGRI observations and even subsequent FY-4B AGRI applications? The algorithm's
427 universal applicability was tested using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

429 Taking the full-disk observation of FY-4A AGRI at 04:00 (UTC, the same below)
430 on 1 June 2023 as an example, the radiance observations from 14 channels are initially
431 fed into the random forest cloud detection model to determine the sky classification
432 (overcast, partly cloudy or clear sky) in each AGRI field. The random forest cloud
433 fraction retrieval model is utilized to retrieve the cloud fraction in scenes identified as
434 partly cloudy. Figure 3(a) is the observed albedo at $0.67\ \mu\text{m}$, where the circles represent
435 the contours of the sunglint angle, (b) is the cloud fraction retrievals from random forest
436 algorithm, (c) is the official operational cloud fraction product and (d) is random forest
437 cloud fraction retrievals with sun-glint correction. It can be seen from Figure 3 that

438 many clear-sky scenes are erroneously identified as cloudy by the operational product
439 and the cloud fraction is generally overestimated with many scenes having a cloud
440 fraction of 1. The random forest algorithm identifies more regions as clear skies or
441 partly cloudy than the operational products, matching better with the observations in
442 the 0.67 μm albedo image. Brighter regions in the visible image correspond to cloud
443 cover areas and darker areas represent clear sky conditions. The sun glint region in the
444 central South China Sea (the circled area in Figure 3(a)) is depicted in Figure 3(b),
445 where the clear-sky scenes over the ocean are misidentified as partly cloudy by random
446 forest algorithm due to the increase in observed albedo. Although operational product
447 in this area also suffers from the impact of unremoved sun glint, it identifies more clear-
448 sky scenes and the cloud fraction is relatively low. Thus, it is evident that the random
449 forest algorithm exhibits significant cloud detection and cloud fraction errors in these
450 sun glint regions. Correction is necessary for the cloud fraction retrievals in the sun
451 glint region.

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

458



459

460

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

465 Statistical analysis was conducted on the correction effect using samples with sun
 466 glint in the training data. The POD and FAR in sun glint area is listed in table 5 and the
 467 error is in table 6. It can be seen that after correcting for cloud fraction, the POD for
 468 clear skies has increased from 0.0987 to 0.9023. The FAR for partly cloudy has

469 decreased from 0.7943 to 0.0276. Both ME, MAE, and RMSE show significant
 470 reductions, and the results after correction outperform operational products.

471 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

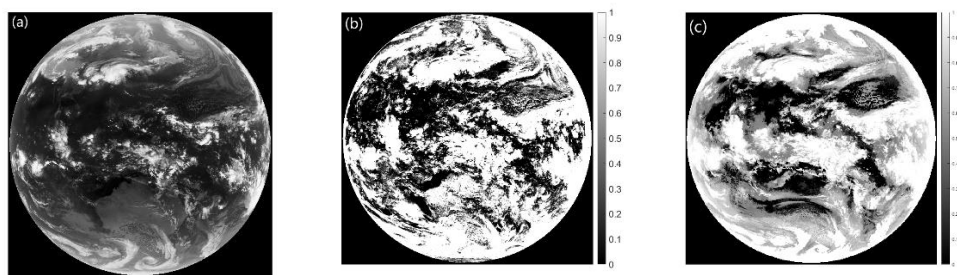
472

473 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

474 FY-4B launched in 2021 has a total of 15 channels with an additional low-level
 475 water vapor channel at 7.42 μm compared to FY-4A. Taking the full-disk observation
 476 of FY-4B AGRI at 17:00 on April 18, 2023, as an example, The radiance observation

477 data of the remaining eight channels (near-infrared and infrared channels) except for
478 the 7.42 μm channel and the visible light channels were input into the random forest
479 cloud detection model. Figure 4 (a) shows the brightness temperature distribution
480 observed in the 10.8 μm channel of FY-4B AGRI, (b) represents the operational cloud
481 fraction product for FY-4B AGRI and (c) shows the cloud fraction retrieved by this
482 algorithm. Figure 4 illustrates that the random forest algorithm identifies more regions
483 as clear skies or partly cloudy than the operational products, aligning better with the
484 brightness temperature observations in 10.8 μm . Especially in high latitude regions of
485 the southern hemisphere and areas with strong convection near the equator, the cloud
486 cover provided by operational products is too high and even misjudged. It can be seen
487 that the random forest algorithm is also suitable for cloud fraction retrieval of FY-4B
488 AGRI.



489
490 **Figure 4** FY-4B AGRI at 17:00 on 18 April 2023, (a) brightness temperature of
491 10.8 μm channel, (b) operational cloud fraction product, (c) random forest cloud
492 fraction retrieval.

493

494 **4 Conclusion**

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

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

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

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

513 (4) The sun-glint area cloud fraction correction curve, fitted with SunGlintAngle

514 as the weight, greatly improves the accuracy of cloud fraction retrieval and reduces the
515 misclassification rate of clear sky scenes as partly cloudy or partly cloudy scenes as
516 overcast due to increased reflectance.

517

518 *Data availability*

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

521

522 *Author contributions*

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

526

527 *Competing interests*

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

529

530 *Disclaimer*

531 *Acknowledgements*

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

534 *We acknowledge the High Performance Computing Center of Nanjing University of*

535 *Information Science & Technology for their support of this work.*

536 **References**

- 537 Amato, U., Antoniadis, A., Cuomo, V., Cutillo, L., Franzese, M., Murino, L., Serio,
538 C.: Statistical cloud detection from SEVIRI multispectral images, *Remote Sensing*
539 *of Environment*, 112, 750–766, <https://doi.org/10.1016/j.rse.2007.06.004>, 2008.
- 540 Baum, B., Trepte Q.: A Grouped Threshold Approach for Scene Identification in
541 AVHRR Imagery, *Journal of Atmospheric & Oceanic Technology*, 16, 793-800,
542 [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.
- 543 Breiman L.1999. Random Forests-Random Features [J]. *Machine Learning*.45(1): 5-32.
- 544 Breiman, L. Random Forests. *Machine Learning* **45**, 5–32 (2001).
545 doi.org/10.1023/A:1010933404324
- 546 Merchant, C.J., Harris, A.R., Maturi, E., Maccallum S.: Probabilistic physically based
547 cloud screening of satellite infrared imagery for operational sea surface temperature
548 retrieval, *Quarterly Journal of the Royal Meteorological Society*, 131, 2735-2755,
549 <https://doi.org/10.1256/qj.05.15>, 2005.
- 550 Gao, J., Jing, Y.: Satellite Remote Sensing Cloud Detection Method Based on Fully
551 Convolutional Neural Network, *Infrared Technology*, 41, 607-615, 2019.
- 552 Gomez-Chova, L., Camps-Valls, G., Amoros-Lopez, J., Guanter, L., Alonso, L.,
553 Calpe, J., Moreno, J.: New Cloud Detection Algorithm for Multispectral and
554 Hyperspectral Images: Application to ENVISAT/MERIS and PROBA/CHRIS

555 Sensors, *IEEE International Symposium on Geoscience and Remote Sensing*, 2757–
556 2760, doi:10.1109/igarss.2006.709, 2006.

557 Hu, J.: Research on Cloud Detection Algorithm of Remote Sensing Image Based on
558 Convolution Neural Network, *Nanjing University of Information Science and
559 Technology*. doi:10.27248/d.cnki.gnjqc.2020.000625, 2020.

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

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

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

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

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

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

576 10.6046/zrzyyg.2021259. 2022

577 Rossow, W. B., Leonid, C.G.: Cloud detection using satellite measurements of
578 infrared and visible radiances for ISCCP. *Journal of Climate*, 12, 2341-2369,
579 [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.

580 R.O. Duda, P.E. Hart, and D.G. Stork, Pattern Classification, New York: John Wiley
581 & Sons, 2001, pp. xx + 654, ISBN: 0-471-05669-3. *Journal of Classification* 24,
582 305–307 (2007). <https://doi.org/10.1007/s00357-007-0015-9>

583 Solvsteen, C.: Correlation based cloud-detection and an examination of the split-
584 window method, *Proceedings of SPIE - The International Society for Optical*
585 *Engineering*, 86-97, 1995.

586 Wang, Z.: CloudSat Project: CloudSat 2B-CLDCLASS-LIDAR product process
587 description and interface control document, *Jet Propulsion Laboratory*, 2019.

588 Yan, J., Guo, X., Qu, J.: An FY-4A/AGRI cloud detection model based on the naive
589 Bayes algorithm, *Remote Sensing for Natural Resources*, 34, 33-42, 2022.

590 Zhang, W., He, M., Mak, M.W.: Cloud detection using probabilistic neural networks,
591 *Geoscience and Remote Sensing Symposium*, IEEE 2373-2375, 2001.

592 Zhang, Y., William, B. R., Andrew, A. L., Valdar, O., Michael, I. M.: Calculation of
593 radiative fluxes from the surface to the top of atmo- sphere based on ISCCP and
594 other global data sets: Refine- ments of the radiative transfer model and the input
595 data, *Journal of Geophysical Research Atmospheres*, 109, 1-27,
596 <https://doi.org/10.1029/2003JD004457>, 2004.