1 **Using machine learning algorithm to retrieve cloud fraction based on**

2 **FY-4A AGRI observations**

3 Jinyi Xia[1]  Li Guan[1]

4 [1]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

8

9 **Abstract**

10 Cloud fraction as a vital component of meteorological satellite products plays an

11 essential role in environmental monitoring, disaster detection, climate analysis and

12 other research areas. A long short-term memory (LSTM) machine learning algorithm

13 is used in this paper to retrieve the cloud fraction of AGRI (Advanced Geosynchronous

14 Radiation Imager) onboard FY-4A satellite based on its full-disc level-1 radiance

15 observation. Correction has been made subsequently to the retrieved cloud fraction in

16 areas where solar glint occurs using a correction curve fitted with sun-glint angle as

17 weight. The algorithm includes two steps: the cloud detection is conducted firstly for

18 each AGRI field of view to identify whether it is clear sky, partial cloud or overcast

19 cloud coverage within the observation field. Then the cloud fraction is retrieved for the

20 scene identified as partly cloudy. The 2B-CLDCLASS-LIDAR cloud fraction product

21 from Cloudsat& CALIPSO active remote sensing satellite is employed as the truth to

22    assess the accuracy of the retrieval algorithm. Comparison with the operational AGRI

23    level 2 cloud fraction product is also conducted at the same time. During daytime, the

24    probability of detection (POD) for clear sky, partly cloudy, and overcast scenes in the

25    official operational cloud detection product were 0.5359, 0.7041, and 0.7826,

26    respectively. The POD for cloud detection using the LSTM algorithm were 0.8294,

27    0.7223, and 0.8435. While the operational product often misclassified clear sky scenes

28    as cloudy, the LSTM algorithm improved the discrimination of clear sky scenes, albeit

29    with a higher false alarm rate compared to the operational product. For partly cloudy

30    scenes, the mean error (ME) and root-mean-square error (RMSE) of the operational

31    product were 0.2374 and 0.3269. The LSTM algorithm exhibited lower ME (0.1134)

32    and RMSE (0.1897) than the operational product. The large reflectance in the sun-glint

33    region resulted in significant cloud fraction retrieval errors using the LSTM algorithm.

34    However, after applying the correction, the accuracy of cloud cover retrieval in this

35    region greatly improved. During nighttime, the LSTM model demonstrated improved

36    POD for clear sky and partly cloudy scenes compared to the operational product, while

37    maintaining a similar POD value for overcast scenes and a lower false alarm rate. For

38    partly cloudy scenes at night, the operational product exhibited a positive mean error,

39    indicating an overestimation of cloud cover, whereas the LSTM model showed a

40    negative mean error, indicating an underestimation of cloud cover. The LSTM model

41    also exhibited a lower RMSE compared to the operational product.

42    **Key words:** Cloud detection, cloud fraction, FY-4A AGRI, LSTM neural network.

**Introduction**

43

44      Clouds occupy a significant proportion within satellite remote sensing data

45      acquired for Earth observation. According to the statistics from the International

46      Satellite Cloud Climatology Project (ISCCP), the annual average global cloud coverage

47      within satellite remote sensing data is around 66% with even higher cloud coverage in

48      specific regions (such as the tropics) (Zhang , et al., 2004). The impact of clouds on the

49      radiation balance of the Earth's atmospheric system is determined by the optical

50      properties of clouds. Cloud detection, as a vital component of remote sensing image

51      data processing, is considered a critical step for the subsequent identification, analysis,

52      and interpretation of remote sensing images. Therefore, accurately determining cloud

53      coverage is essential in various research domains, such as environmental monitoring,

54      disaster surveillance and climate analysis.

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

56      launched by China in 2016. The uploaded AGRI (Advanced Geosynchronous Radiation

57      Imager) has 14 channels and captures full-disk observation every 15 minutes. In

58      addition to observing clouds, water vapor, vegetation and the Earth's surface, it also

59      possesses the capability to capture aerosols and snow. Moreover, it can clearly

60      distinguish different phases and particle size of clouds and obtain high- to mid-level

61      water vapor content. It is particularly suitable for cloud detection due to its

62      simultaneous use of visible, near-infrared and long-wave infrared channels for

63    observation with high spatial resolution.

64         Numerous cloud detection algorithms have been provided based on observations

65    from satellite-borne imagers. The threshold method has been widely employed by

66    researchers, encompassing the early ISCCP (International Satellite Cloud Climatology

67    Project) method (Rossow, 1993) and the proposed threshold methods based on different

68    spectral features or underlying surfaces. Kegelmeyer (1994) used a straightforward

69    cloud pixel as threshold for cloud detection with Whole Sky Imaging Cameras.

70    Solvsteen (1995) distinguished cold water pixels and cloud pixels by analyzing the

71    correlation between different channels based on AVHRR (Advanced Very High

72    Resolution Radiometer) images. A grouping threshold method based on AVHRR

73    images has been developed by Baum and Trepte (1996) to classify scenes as clouds,

74    fires, smoke or snow. LI and Zhang (2006) proposed a multispectral integrated cloud

75    detection algorithm based on the characteristics of MODIS instrument channels and the

76    spectral characteristics of different objects (clouds, snow, land, etc.). Zhang et al. (2020)

77    used a multi-temporal cloud detection method based on FY-4A AGRI data to identify

78    observations on the Qinghai-Tibet Plateau. However, there is a significant subjectivity

79    in selection of thresholds whether it is the single and fixed threshold in the early days,

80    multiple thresholds, dynamic thresholds, or adaptive thresholds. These thresholds are

81    highly influenced by factors such as season and climate.

82         The other category of cloud detection algorithms is the based on statistical

83    probability theory. Such as the principal component discriminant analysis and quadratic

84    discriminant analysis methods were used to SEVIRI (Spinning Enhanced Visible and

85    Infrared Imager) cloud detection (Amato et al., 2008). The cloud detection algorithm

86    for Thermal Infrared (TIR) sensor was based on the Bayesian theory of total probability

87    (Merchant et al., 2010) and the naive Bayes algorithm for AGRI (Qu , et al., 2022). The

88    unsupervised clustering cloud detection algorithms for MERIS (Medium Resolution

89    Imaging Spectrometer) (GomezChova , et al., 2007) and the fuzzy C-means clustering

90    algorithms for MODIS (Pan, et al., 2009) all have achieved high accuracy in cloud

91    detection.

92        More and more machine learning algorithms are being utilized by researchers in

93    cloud detection studies with the development of machine learning. For instance, the

94    probabilistic neural networks, especially radial basis function networks was used for

95    AVHRR cloud detection (Zhang, et al., 2001). The utilization of convolutional neural

96    network methods (Hu, et al., 2020) offers important perspectives for cloud detection

97    research.

98        Currently, there is limited research literature on cloud detection and cloud fraction

99    retrieval algorithms for FY-4A/4B AGRI. The operational cloud fraction product of FY-

100   4A AGRI utilized a threshold method with 4 km spatial resolution. Differences in

101   climatic and environmental factors lead to varying albedo and brightness temperature

102   observations for the instrument at different times and locations. Therefore, the choice

103   of thresholds is easily influenced by factors such as season, latitude and land surface

104   type (Gao and Jing, 2019). Using multiple sets of thresholds for discrimination would

105     significantly slow down the cloud detection process. Moreover, most algorithms focus

106     solely on cloud detection, which classified the observed scenes into cloud or clear-sky

107     without providing the specific cloud fraction information for the scenes.

108         In summary, a LSTM (Long Short-Term Memory) machine learning algorithm for

109     cloud fraction retrieval was established using level-1 radiation observations from FY-

110     4A AGRI full-disk scanning in this paper. The cloud fraction of the level-2 product 2B-

111     CLDCLASS-LIDAR from Cloudsat&CALIPSO was used as the reference label. The

112     retrievals were compared against with the cloud fraction of 2B-CLDCLASS-LIDAR

113     and the AGRI operational products to verify the algorithm accuracy.

114     **1 Research Data and Preprocessing**

115     *1.1   FY-4A data*

116     FY-4A was successfully launched on December 11, 2016. Starting from May 25, 2017,

117     FY-4A drifted to a position near the main business location of the Fengyun

118     geostationary satellite at 104.7 degrees east longitude on the equator. Its successful

119     launch marked the beginning of a new era for China's next-generation geostationary

120     meteorological satellites as an advanced comprehensive atmospheric observation

121     satellite. The Advanced Geosynchronous Radiation Imager (AGRI), one of the main

122     payloads of the Fengyun-4 series geostationary meteorological satellites, can perform

123     large-disk scans and rapid regional scans at a minute level. It has total 14 observation

124     channels with the main task of acquiring cloud images. The channel parameters and

125     main uses of AGRI are detailed in Table 1. FY-4A AGRI data was downloaded from

126     the official website of the China national satellite meteorological center

127     (http://satellite.nsmc.org.cn), including level-1 full disk radiation observation data

128     preprocessed through quality control, geolocation and radiation calibration as well as

129     level-2 cloud fraction product (CFR). The spatial resolution of these data is all 4 km

130     and the temporal resolution is 15 minutes.

131                          **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 | 11.5 0~ 12.50 | 12.0 | 4 | surface and cloud-top temperatures |
| 14 | 13.2 ~ 13.8 | 13.5 | 4 | cloud-top height |

132

133     **1.2  CloudSat & Calipso Cloud Product**

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

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

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

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

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

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

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

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

142    CloudSat is also among NASA's A-Train series satellites. The CPR (Cloud Profile

143    Radar) installed on it operates at 94 GHz millimeter-wave and is capable of detecting

144    the vertical structure of clouds and providing vertical profiles of cloud parameters. The

145    scanning wavelengths of CPR and CALIOP are different. CALIOP is capable of

146    observing the top of mid-to-high level clouds, whereas CPR can penetrate optically

147    thick clouds. Combining the strengths of these two instruments enables the acquisition

148    of precise and detailed information on cloud layers and cloud fraction.

149        The joint level 2 product 2B-CLDCLASS-LIDAR is mainly utilizing in this study.

150    It provides the cloud fraction at different heights with horizontal resolution 2.5 km

151    (along-track) × 1.4 km (cross-track) through combining the observations from CPR and

152    CALIOP (Zhen, et al., 2018). The CloudSat product manual (Wang, 2019) can be

153    referred for more detailed information on 2B-CLDCLASS-LIDAR. The data used is

154    available    for    download    from    the    ICARE    data    and    services    center

155    (https://www.icare.univ-lille.fr/data-access/data-archive-access/).

### 1.3 Establishment of Training Data

156

157 The crucial aspect of establishing a training data in machine learning algorithms

158 is how to obtain the cloud fraction values (ground truth) as labels. The error in cloud

159 fraction retrieved solely from passive remote sensing instruments is significant. Using

160 active remote sensing data can provide more accurate cloud fraction information in the

161 vertical direction. Therefore, the spatiotemporally matched 2B-CLDCLASS-LIDAR

162 cloud fraction are utilized as output labels in this paper.

163 The FY-4A AGRI and 2B-CLDCLASS-LIDAR data with a distance difference

164 between fields of view within 1.5 km and a time difference within 15 minutes are

165 spatiotemporal matched. To make the 2B-CLDCLASS-LIDAR cloud fraction data

166 collocated within AGRI pixels more effective, at least two 2B-CLDCLASS-LIDAR

167 pixels are required within each AGRI field of view. The cloud fraction average of these

168 pixels is used as the cloud fraction for that AGRI pixel.

169 Cloud detection and cloud fraction label generation for 2B-CLDCLASS-LIDAR

170 are as follows. There may be multiple layers of clouds in each field of view. If there is

171 at least one layer cloud with cloud fraction of 1 in the 2B-CLDCLASS-LIDAR profile,

172 then the scene is labeled as overcast with a cloud fraction of 1. If all layers in the profile

173 are cloud-free, the scene is labeled as clear sky. The scene between the above two

174 situations is labeled as partly cloudy and the cloud fraction is the average of cloud

175 fractions at different layers.

176    The algorithm includes two steps: the cloud detection is conducted firstly for each

177    AGRI field of view to identify whether it is clear sky, partial cloud or overcast cloud

178    coverage within the observation field. Then the cloud fraction is retrieved for the scene

179    identified as partly cloudy. So the training data include A dataset used for cloud

180    detection and B dataset for cloud fraction retrieval.    The input variables in A dataset

181    are the FY-4A AGRI level-1 radiative observations from 14 channels and the output

182    variable is the temporally and spatially matched 2B-CLDCLASS-LIDAR cloud

183    detection label. The output is categorized into three types: overcast, partly cloudy and

184    clear sky with values 1, 2 and 3 respectively. To ensure diversity and representativeness

185    of the samples, the three conditions of overcast, partly cloudy, and clear sky each

186    account for one-third of the sample size in dataset A. Regarding the samples for partly

187    cloudy type in dataset A, the collocated 2B-CLDCLASS-LIDAR cloud fraction

188    products serve as output labels for cloud fraction retrieval model B. The input of

189    training dataset B remains the FY-4A AGRI level-1 radiative observations.

190    Due to the lifespan of the instrument only 2B-CLDCLASS-LIDAR data before

191    July 2019 can be obtained. So, the FY-4A AGRI observations and 2B-CLDLASS-

192    LIDAR matched in time and space in May 2019 are used as training samples to build

193    the algorithm model. The paired samples of whole June 2019 are served as the testing

194    samples to assess the model's retrieval accuracy. The number of training samples in

195    May are 12,420 for dataset A and 4140 for B. Testing samples in June are 15,459 for A

196    and 5,153 for B.

197     Although the retrieval model was trained and tested using 2019 data, the algorithm

198     was also applied to real-time observations of FY-4A and FY-4B AGRI in 2023 to verify

199     its universality.

200

201     **2.  Long Short-Term Memory (LSTM) Algorithm**

202         LSTM is an improved algorithm based on RNN (Recurrent Neural Network) with

203     the ability to retain long-term memory. and demonstrates improved performance in

204     longer sequences data comparing to ordinary RNNs (Sarker, 2001).    It can effectively

205     address the challenges of gradient explosion and gradient vanishing over time in

206     models., LSTM network has been extensively applied in diverse domains owing to its

207     distinctive features, such as meteorology and environmental prediction and so on (Bao,

208     et al., 2024; Bai and Shen. 2019). The structure of the LSTM unit is depicted in Figure

209     1. The update and transmission of historical information is facilitated through the

210     internal control of three states: the Forget Gate, the Input Gate and the Output Gate.

211     The pertinent mathematical expressions are:

212         $f_t = \sigma(W_f^T \times [h_{t-1}, x_t] + b_f)$                                    (1)

213         where $f_t$ denotes the output of the Forget Gate, $\sigma$ signifies the Sigmoid

214     activation function; $W_f^T$ and $b_f$ correspond to the weight and bias of the Forget Gate,

215     respectively, $x_t$ stands for the current input, $h_{t-1}$ represents the output from the

216     previous time step.

217 $$i_t = \sigma(W_i^T \times [h_{t-1}, x_t] + b_i) \tag{2}$$

218     where $i_t$ represents the information updated after $\sigma$ activation, $W_i^T$ and $b_i$

219     denote the weight and bias, respectively.

220 $$\widehat{C_t} = \sigma(W_c^T \times [h_{t-1}, x_t] + b_c) \tag{3}$$

221     $\widehat{C_t}$ signifies the information updated after tanh activation, $W_c^T$ and $b_c$ denote

222     the weight and bias, respectively.

223 $$C_t = f_t \times C_{t-1} + i_t \times \widehat{C_t} \tag{4}$$

224     $C_t$ is the current information of the LSTM structure, $C_{t-1}$ denotes the

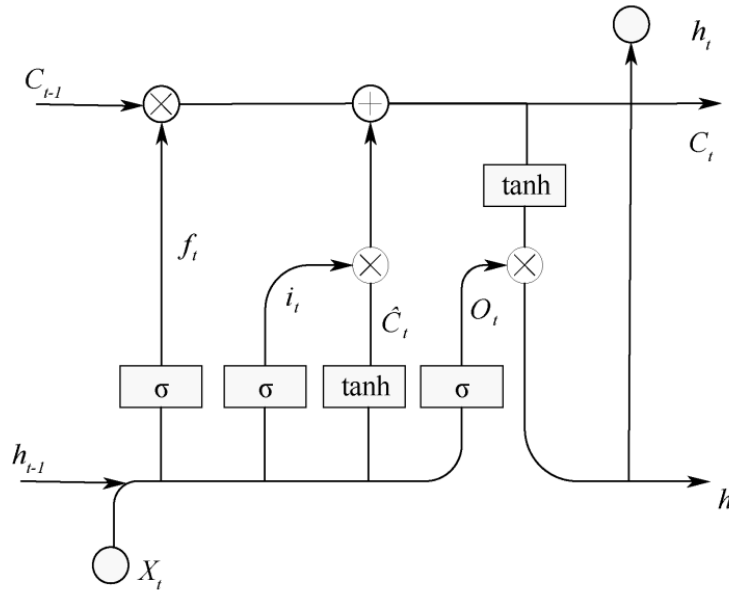225     information of the LSTM structure from the previous time step.

226 $$O_t = \sigma(W_O^T \times [h_{t-1}, x_t] + b_O) \tag{5}$$

227     $O_t$ is the current output information, $W_O^T$ and $b_O$ denote the weight and bias,

228     respectively.

229 $$h_t = o_t \times \tanh(C_t) \tag{6}$$

230     $h_t$ denotes the current output result.

**Figure 1** LSTM cell structure (Kong, et al., 2018)

231

232

233      In a neural network, the hidden layer is a layer or multiple layers located between

234    the input layer and the output layer. Each hidden layer consists of multiple nodes, which

235    process the input data and generate outputs through connection weights and activation

236    functions. Increasing the size of the hidden layer can enhance the network's

237    representational capacity and learning ability, as more nodes can capture additional data

238    patterns and features. However, having a hidden layer that is too large may lead to

239    overfitting, making the network overly complex and difficult to train. Typically, the

240    optimal size of the hidden layer is determined by trying different sizes and evaluating

241    their performance on a validation set. The hidden layer sizes for both the cloud

242    classification model and the cloud fraction retrieval model in this paper are set to 3.

243      The key model parameter 'batch size' has two main impacts on training network:

244    (1) A larger batch size typically reduces the training time per epoch as more samples

245    are processed with each parameter update. On the contrary, a smaller batch size may

246    slow down the training speed since more iterations are needed to complete an epoch.

247    (2) Model Performance: Different batch sizes can impact the model performance.

248    Generally, a larger batch size may lead to quicker model convergence, yet it could

249    increase the risk of overfitting at times; whereas a smaller batch size could aid in the

250    model's generalization ability but might result in a less stable training process. In this

251    paper, the batch size of the model is set to 500.The optimizer is configured with the

252    Adam gradient descent algorithm, and the loss function used is cross-entropy.

253        The training dataset A was used to construct the LSTM cloud detection model. For

254    daytime, the inputs are the radiation observations from 14 channels of FY-4A AGRI

255    with 'input size' 14. However, during nighttime, as there are no observations in the

256    visible light channels (channels 1 to 6) of AGRI, the inputs consisted of the radiance

257    observations of channels 7 to 14 of FY-4A AGRI with 'input size' 8. The output label

258    is the classification of field of view, including overcast, partly cloudy and clear sky.

259        To derive the specific cloud fraction for AGRI scenes identified as partly cloudy

260    in the previous cloud mask step, an LSTM cloud fraction retrieval model needs to be

261    constructed. The training dataset B was used to train the cloud fraction retrieval model.

262    For daytime, the input is the observed radiances for all channels of AGRI (input

263    size=14), while during nighttime, the input comprises the observed radiance values of

264    channels 7 to 14 of AGRI (input size = 8). The output label is the value of cloud fraction

265    in the scene ranging from 0 to 1. When selecting parameters for the LSTM cloud

266    fraction model, a batch size of 60 is chosen due to the limited sample number in dataset

267    B. The optimizer is also configured with the Adam gradient descent algorithm. The loss

268    function used is mean square error.

269    **3.   Results and Analysis**

270    To assess the accuracy and stability of the retrieval model, two types of validation

271    methods are utilized. One way involves a direct comparison from images, qualitatively

272    comparing the model's retrieval results and official cloud fraction products with AGRI

273    observed cloud images. Another way is quantitative comparison using 2B-

274    CLDCLASS-LIDAR as the true value. Four quantitative parameters, including

275    possibility of detection(POD), alse alarm rate(FAR), mean error (ME) and root mean

276    square error (RMSE) are introduced. 'Possibility of detection' is calculated using the

277    formula POD=TP/(TP+FN), and false alarm rate is calculated using the formula

278    FAR=FP/(TP+FP). Taking the covercast scenes as an example, TP represents the

279    number of correctly identified overcast, FN represents the number of overcast scenes

280    wrongly identified as partly cloudy or clear sky, and FP represents the number of clear

281    sky or partly cloudy scenes wrongly identified as overcast.The ME (mean error) and

282    RMSE (root mean square error) are utilized to assess the accuracy of the LSTM cloud

283    fraction model in retrieving cloud fraction for partly cloudy scenes.

### 3.1 Objective Analysis of Cloud Fraction Retrievals

284

285     The test samples from dataset A (i.e., June data) are used to perform cloud

286 detection experiments based on the cloud detection model mentioned above. The

287 temporally and spatially matched 2B CLDCLASS-LIDAR cloud mask products are

288 used as reference to evaluate the accuracy of cloud detection. The POD and FAR for

289 different view field classifications are shown in Table 2. Columns 2 and 4 represent the

290 operational cloud detection products for daytime and nighttime respectively, for the

291 same time and pixel. Columns 3 and 5 represent the LSTM cloud detection results for

292 daytime and nighttime respectively. The table indicates that during daytime, operational

293 cloud detection products have a relatively low possibility of detection for clear sky view

294 fields. However, the LSTM model increases the possibility of detection for clear sky

295 from 0.54 to 0.83. Moreover, for some partly cloudy and overcast view fields, the

296 possibilities of detection is higher than those of operational cloud detection products.

297 During nighttime, compared to operational cloud detection products, the LSTM model

298 increases the POD for clear sky from 0.51 to 0.73, with slightly higher possibilities of

299 detection for partial cloud view fields than the operational products, while the

300 possibility of detection for full cloud view fields is lower. During the day, the

301 Operational product has a lower false alarm rate for clear sky compared to the LSTM

302 model, while the LSTM model has a lower false alarm rate for partly cloudy and

303 overcast conditions than the Operational product. At night, the LSTM model

304   significantly reduces the false alarm rate for overcast conditions compared to the

305   Operational product.

306                    **Table 2** POD and FAR of Cloud Detection

|  | Sky Classification | Daytime Operational Cloud Detection Product | Daytime LSTM Results | Nighttime Operational Cloud Detection Product | Nighttime LSTM Results |
|---|---|---|---|---|---|
| POD | Clear Sky | 0. 5359 | 0.8294 | 0.5136 | 0.7341 |
|  | Partly cloudy | 0.7041 | 0.7223 | 0.6957 | 0.7101 |
|  | Overcast | 0.7826 | 0.8435 | 0.7984 | 0.7523 |
| FAR | Clear Sky | 0.2174 | 0.3633 | 0.1789 | 0.1983 |
|  | Partly cloudy | 0.2959 | 0.1677 | 0.3107 | 0.3488 |
|  | Overcast | 0.4641 | 0.2358 | 0.5543 | 0.2105 |

307

308   For the view fields judged as partly cloudy by the aforementioned model, the cloud

309   amount in the AGRI view field was inverted using the LSTM cloud amount model

310   established earlier in this text. For samples classified as partly cloudy by the model,

311   operational products and 2B-CLDCLASS-LIDAR cloud amount products, the mean

312   error and root mean square error (RMSE) of the cloud amount retrieval were calculated

313   based on the matched 2B-CLDCLASS-LIDAR cloud amount product as ground truth,

314   separately for daytime and nighttime operational cloud amount products (columns 2

315   and 4) and the LSTM-inverted cloud amount (columns 3 and 5), as shown in Table 3.

316   It can be observed that during daytime, compared to the FY-4A operational cloud

317   amount product, the LSTM cloud amount retrieval model shows significant

318   improvement in both mean error (ME) and RMSE. The ME decreases from 0.23 to 0.11,

319 and the RMSE decreases from 0.32 to 0.19, indicating that the LSTM cloud amount

320 retrieval model provides more accurate estimates of cloud amount. For nighttime, the

321 ME of the operational cloud amount product is positive, indicating an overall

322 overestimation of cloud amount. In contrast, the ME of the LSTM model is negative,

323 indicating an overall underestimation of cloud amount. The RMSE of the LSTM model

324 retrieval results during nighttime is lower than that of the operational cloud amount

325 product.

326

**Table 3** Errors in cloud fraction retrieval

|  | Daytime Operational Cloud Detection Product | Daytime LSTM Results | Nighttime Operational Cloud Detection Product | Nighttime LSTM Results |
|---|---|---|---|---|
| ME | 0.2374 | 0.1134 | 0.2488 | -0.1911 |
| RMSE | 0.3269 | 0.1897 | 0.3374 | 0.2361 |

327 **3.2 Cloud fraction correction in sun glint regions**

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

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

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

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

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

333 increase in visible channel observation albedo will affect various subsequent

334 applications of data, including cloud detection and cloud cover retrieval, etc.

335    The position of Sun glint area can be determined using the SunGlintAngle value

336    in the FY-4A GEO file. SunGlintAngle is defined as the angle between the satellite

337    observation direction or reflected radiation direction and the mirror reflection direction

338    on a calm surface (horizontal plane). It is generally accepted that the range of

339    SunGlintAngle < 15° is easily affected by sun glint (Kay S, et al., 2009). The positions

340    of the SunGlintAngle contour lines at 5 and 15° are marked in Figure 2(a). It can be

341    observed that the edge of sun glint in Figure 2(a) essentially overlaps with the position

342    of SunGlintAngle = 15°. Thus, the region where SunGlintAngle < 15° is defined as the

343    sun glint range in this paper and only the cloud fraction within this range will be

344    adjusted in the subsequent correction.

345    To correct the cloud fraction in the sun glint region, we initially identified 672

346    fields of view where sun glint occurred in the FY-4A AGRI observations between 1

347    June and 31 July 2019.  Subsequently, a direct least squares fitting was conducted

348    between the inverted cloud fraction and the collocated 2B-CLDCLASS-LIDAR cloud

349    fraction (ground truth). The scatter plot is illustrated in Figure 2(b), where x-axis is the

350    2B-CLDCLASS-LIDAR cloud fraction and y-axis is the model-inverted cloud fraction.

351    The blue line represents the curve (namely Eq.7) fitted by the least squares method

352    between the retrievals and the truths. The thin dash line is the x=y line. It is evident that

353    the inverted cloud fraction is generally slightly overestimated.

354    Taking observations at 04:00 on 5 June 2019 as an example, Figure 2(c) presents

355    the distribution of SunGlintAngle and the flight trajectory of the Cloudsat&Calypso

356    satellite. White circles denote the sun glint region with SunGlintAngle < 15° and the

357    white line represents the satellite flight track. As depicted in the figure, the majority of

358    Cloudsat&Calypso flight trajectories do not pass through the central position of sun

359    glint area but instead traverse locations with larger SunGliantAngle values. The

360    intensity of sun glint effect decreases with the increase of SunGliantAngle. This

361    suggests that the true values for spatial and temporal matching mostly do not fall within

362    the strongest sun glint region. From Figure 2(d), it can be seen that the impact of sun

363    glint becomes stronger as SunGlintAngle decreasing, which results in a higher

364    observation albedo. This further leads to the overestimated cloud fraction values in the

365    retrieval. It is evident that the cloud fraction error is related to the value of

366    SunGlintAngle and this influence is not considered in Eq. (7). Directly applying

367    equation (7) to correct the cloud fraction retrievals would result in a too small correction

368    intensity for the FOVs near the center of sun glint and an excessively large correction

369    intensity for the FOVs in the Sun-glint edge region (even erroneous clear sky may

370    appear). Considering this, a correction formula (8)-(9) using SunGlintAngle as weight

371    is introduced, where $W_i$ represents the angle weight for a certain pixel $i$ in the sun glint

372    region, n is the number of pixels within the SunGlintAngle < 15° range, yi is the initial

373    model retrieval of cloud cover for the field of view $i$ and $x_i$ is the final corrected cloud

374    fraction.

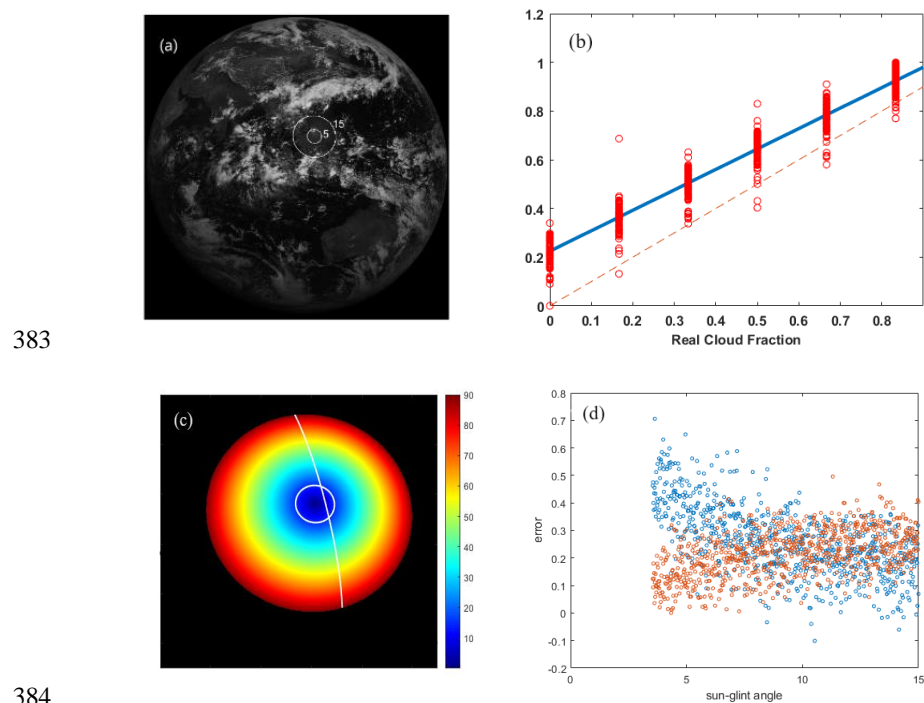375    $$x = (y - 0.2562)/0.8428 \tag{7}$$

376
$$W_i = \frac{glintangle_i}{\frac{1}{n}\Sigma_{i=0}^n glintangle_i} \tag{8}$$

377
$$x_i = W_i \left(\frac{y_i - 0.2526}{0.8428}\right) \tag{9}$$

378 Figure 2(d) shows the distribution of errors with respect to SunGlintAngle,

379 where the blue dots represent the error distribution corrected using formula

380 (7), and the orange dots represent the error distribution corrected using

381 formula (9). It can be seen from Figure 2(d) that after correction by formula

382 (9), the errors in the smaller range of SunGlintAngle are significantly reduced.

383



384

385 **Figure 2** (a) albedo image of 0.67μm channel (the circles are the contours of the sun-

386 glint angle), (b) Scatter plot of cloud fraction in sun glint region, (c) Distribution of

387 SunGlintAngle and satellite flight track of CloudSat & Calypso at 4:00 on June 5, 2019,

388    (d) Distribution of cloud fraction retrieval error with sun-glint angle.

389    **3.3 Algorithm universal applicability testing**

390          Although the retrieval model in this article was built based on data from May 2019

391    due to the limited lifespan of the instrument, how effective is it in real-time FY-4A

392    AGRI observations and even subsequent FY-4B AGRI applications? The algorithm's

393    universal applicability was tested using real-time observations from FY-4A and FY-4B
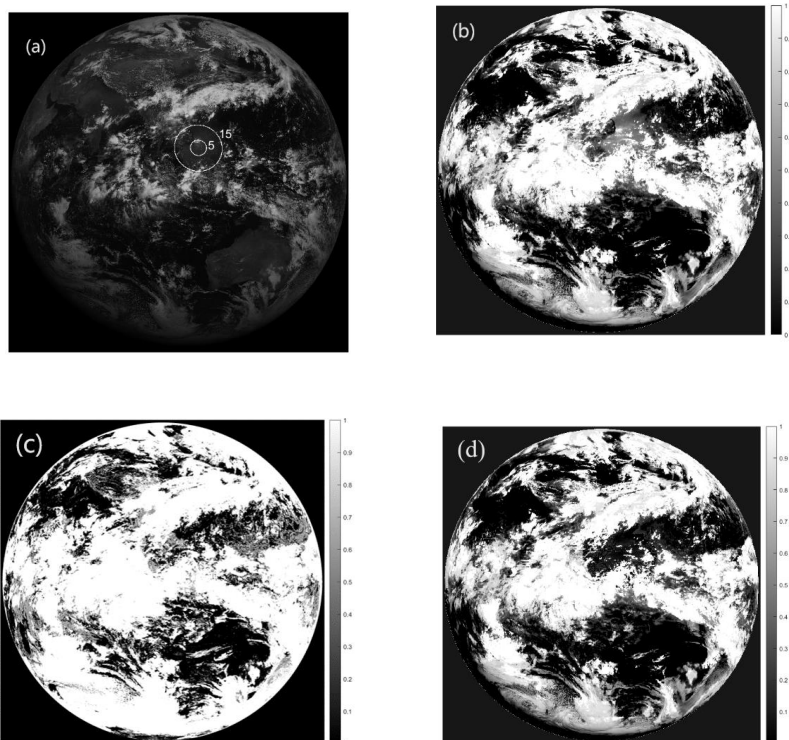
394    AGRI in 2023.

395          Taking the full-disk observation of FY-4A AGRI at 04:00 (UTC, the same below)

396    on 1 June 2023 as an example, the radiance observations from 14 channels are initially

397    fed into the LSTM cloud detection model to determine the sky classification (overcast,

398    partly cloudy or clear sky) in each AGRI field. The LSTM cloud fraction retrieval

399    model is utilized to estimate the cloud fraction in scenes identified as partly cloudy.

400    Figure 3(a) is the observed albedo at 0.67 μm, where the circles represent the contours

401    of the sunglint angle, (b) is the cloud fraction retrievals from LSTM algorithm, (c) is

402    the official operational cloud fraction product and (d) is LSTM cloud fraction retrievals

403    with sun-glint correction. It can be seen from Figure 3 that many clear-sky scenes are

404    erroneously identified as cloudy by the operational product and the cloud fraction is

405    generally overestimated with many scenes having a cloud fraction of 1. The LSTM

406    algorithm identifies more regions as clear skies or partly cloudy than the operational

407    products, matching better with the observations in the 0.67 μm albedo image. Brighter

408 regions in the visible image correspond to cloud cover areas and darker areas represent

409 clear sky conditions. The sun glint region in the central South China Sea (the circled

410 area in Figure 3(a)) is depicted in Figure 3(b), where the clear-sky scenes over the ocean

411 are misidentified as partly cloudy by LSTM algorithm due to the increase in observed

412 albedo. Although operational product in this area also suffers from the impact of

413 unremoved sun glint, it identifies more clear-sky scenes and the cloud fraction is

414 relatively low. Thus, it is evident that the LSTM algorithm exhibits significant cloud

415 detection and cloud fraction errors in these sun glint regions. Correction is necessary

416 for the cloud fraction retrievals in the sun glint region.

417 Figure 3(d) shows the cloud fraction distribution after correction using equation

418 (9) in the sun glint region., The correction eliminates the influence of sun glint

419 comparing to the cloud fraction in sun glint area before correction in Figure 3(b). The

420 scenes misjudged as partly cloudy are corrected to clear sky and match well with the

421 actual albedo observations in 3(a), which accurately restores the true cloud coverage

422 over the South China Sea.

423

424



425

426      **Figure 3** FY-4A AGRI at 04:00 on 1 June 2023 (a) albedo image of 0.67μm channel

427      (the circles are the contours of the sun-glint angle), (b) LSTM cloud fraction

428      retrieval without sun-glint correction, (c) operational cloud fraction product, (d)

429      LSTM cloud fraction retrieval with sun-glint correction.

430      Statistical analysis was conducted on the correction effect using samples with sun

431      glint in the training data. The possibility of detection and false alarm rate in sun glint

432      area is listed in table 4 and the error is in table 5. The possibility of detection for clear

433      skies has increased from 0.09 to 0.83. The false alarm rate for partly cloudy has

434      decreased from 0.89 to 0.17. The mean error of cloud fraction retrievals decreased from

435    0.176 to 0.09. These all indicate that the positive effect of the sun glint correction.

436                    **Table 4** The cloud mask recall rate in sun glint area

|  | Sky Classification | Operational Product | LSTM | LSTM after Correction |
|---|---|---|---|---|
|  | Clear Sky | 0.5535 | 0.0900 | 0.8301 |
| POD | Partly cloudy | 0.6738 | 0.8279 | 0.7436 |
|  | Overcast | 0.8505 | 0.9744 | 0.9744 |
|  | Clear Sky | 0.1437 | 0.0063 | 0.3142 |
| FAR | Partly cloudy | 0.3742 | 0.8972 | 0.1719 |
|  | Overcast | 0.5545 | 0.1324 | 0.1324 |

437

438                    **Table 5** cloud fraction Errors in sun glint area

|  | Operational Product | LSTM Retrievals | LSTM after Correction |
|---|---|---|---|
| ME | 0.2691 | 0.2760 | 0.1634 |
| RMSE | 0.3458 | 0.1948 | 0.1883 |

439        FY-4B launched in 2021 has a total of 15 channels with an additional low-level

440    water vapor channel at 7.42 μm compared to FY-4A. Taking the full-disk observation

441    of FY-4B AGRI at 17:00 on April 18, 2023, as an example, The radiance observation

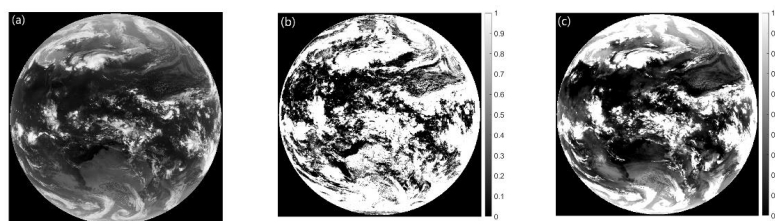442    data of the remaining eight channels (near-infrared and infrared channels) except for

443    the 7.42 μm channel and the visible light channels were input into the LSTM cloud

444    detection model. Figure 4 (a) shows the brightness temperature distribution observed

445 in the 10.8 μm channel of FY-4B AGRI, (b) represents the operational cloud fraction

446 product for FY-4B AGRI and (c) shows the cloud fraction retrieved by this algorithm.

447 Figure 4 illustrates that the LSTM algorithm identifies more regions as clear skies or

448 partly cloudy than the operational products, aligning better with the brightness

449 temperature observations in 10.8 μm. Especially in high latitude regions of the southern

450 hemisphere and areas with strong convection near the equator, the cloud cover provided

451 by operational products is too high and even misjudged. It can be seen that the LSTM

452 algorithm is also suitable for cloud fraction retrieval of FY-4B AGRI.



453

454 **Figure 4** FY-4B AGRI at 17:00 on 18 April 2023, (a) brightness temperature of

455 10.8μm channel, (b) operational cloud fraction product, (c) LSTM cloud fraction

456 retrieval.

457

458 **4 Conclusion**

459 The long short-term memory (LSTM) machine learning algorithm based on FY-

460 4A AGRI full-disc level-1 radiance observations is developed to retrieve the cloud

461     fraction for each field of view in this paper. The accuracy of the algorithm is validated

462     using the 2B CLDCLASS-LIDAR cloud fraction product from the Cloudsat&Calypso

463     active remote sensing satellite and FY-4A AGRI level 2 operational product. The

464     following conclusions are drawn:

465     (1) Not only the cloud detection but also the cloud fraction within each FY-4A

466         AGRI field of view can be retrieved by the LSTM machine learning algorithm.

467     (2) The operational product has a relatively high false alarm rate for clear sky

468         scenes, while the LSTM algorithm improves the probability of detection (POD)

469         for clear sky scenes during the daytime from 0.54 to 0.83. However, the false

470         alarm rate (FAR) is higher compared to the operational product. The POD for

471         clear sky scenes at night increases from 0.51 to 0.73, and the POD for partially

472         cloudy and fully cloudy scenes is comparable to the operational product.

473     (3) For partly cloudy fields, during the day, the mean error and root-mean-square

474         error of the operational product are 0.2374 and 0.3269, respectively, while this

475         algorithm exhibits lower mean error (0.1134) and RMSE (0.1897) than the

476         operational product. At night, the operational product tends to overestimate

477         cloud cover, while this algorithm underestimates cloud cover, with a lower

478         RMSE compared to the operational product.

479     (4) The cloud fraction correction curve for sun glint region fitted with

480         SunGlintAngle as weight significantly improves the accuracy of the LSTM

481         cloud fraction retrievals. It reduces the misjudgment rate where increased

482       albedo leads to the identification of clear-sky scene as partly cloudy or overcast.

483

484 *Data availability*

485 FY-4A AGRI data is available at http://satellite.nsmc.org.cn and the 2B-CLDCLASS-

486 LIDAR data at https://www.icare.univ-lille.fr/data-access/data-archive-access/

487

488 *Author contributions*

489 JX: Formal analysis, Methodology, Software, Visualization and Writing – original draft

490 preparation. LG: Conceptualization, Data curation, Funding acquisition, Supervision,

491 Validation and Writing – review & editing.

492

493 *Competing interests*

494 he contact author has declared that none of the authors has any competing interests.

495

496 *Disclaimer*

497 *Acknowledgements*

500 **References**

501 Bai, S., Shen, X.: PM2. 5 prediction based on LSTM recurrent neural network,

502     *Computer Applications and Software*, 36, 67-70, 2019.

503     Bao S., Qin H., Dai Y.: Short-term precipitation prediction research based on UI-

504     LSTM model, *Radio Engineering*. 1-10, 2023.

505     Baum, B., Trepte Q.: A Grouped Threshold Approach for Scene Identification in

506     AVHRR Imagery, *Journal of Atmospheric & Oceanic Technology*, 16, 793-800,

507     https://doi.org/10.1175/1520-0426(1999)016<0793:AGTAFS>2.0.CO;2, 1999.

508     Merchant, C.J., Harris, A.R., Maturi, E., Maccallum S.: Probabilistic physically based

509     cloud screening of satellite infrared imagery for operational sea surface temperature

510     retrieval, *Quarterly Journal of the Royal Meteorological Society*, 131, 2735-2755,

511     https://doi.org/10.1256/qj.05.15, 2005.

512     Gao, J., Jing, Y.: Satellite Remote Sensing Cloud Detection Method Based on Fully

513     Convolutional Neural Network,*Infrared Technology*, 41, 607-615, 2019.

514     Gomez-Chova, L., Camps-Valls, G., Amoros-Lopez, J., Guanter, L., Alonso, L.,

515     Calpe, J., Moreno, J.: New Cloud Detection Algorithm for Multispectral and

516     Hyperspectral Images: Application to ENVISAT/MERIS and PROBA/CHRIS

517     Sensors*, IEEE International Symposium on Geoscience and Remote Sensing,* 2757–

518     2760, doi:10.1109/igarss.2006.709, 2006.

519     Hu, J.: Research on Cloud Detection Algorithm of Remote Sensing Image Based on

520     Convolution Neural Network, *Nanjing University of Information Science and*

521     *Technology*, doi:10.27248/d.cnki.gnjqc, 2020.

522     Kay, S., Hedley, J., Lavender, S.: Sun Glint Correction of High and Low Spatial

523      Resolution Images of Aquatic Scenes: a Review of Methods for Visible and Near-

524      Infrared Wavelengths, Remote Sensing, 1, 697-730,

525      https://doi.org/10.3390/rs1040697, 2009.

526      Kegelmeyer, W.P.J.: Extraction of cloud statistics from whole sky imaging

527      cameras,1994.

528      Kong, Y.-L., Huang, Q., Wang, C., Chen, J., Chen, J., & He, D. (2018). Long Short-

529      Term Memory Neural Networks for Online Disturbance Detection in Satellite

530      Image Time Series. *Remote Sensing*, 10(3), 452. doi:10.3390/rs10030452

531      Li, W., Zhang, L., Chen, X., Li, D.: The universal cloud detection algorithm of

532      MODIS data, *Society of Photo-Optical Instrumentation Engineers (SPIE)*

533      *Conference Series,* 64190F-64190F-6, doi:10.1117/12.712722 , 2006.

534      Pan, C., Xia B., Chen, Y.: Research on MODIS Cloud Detection Algorithms Based on

535      Fuzzy Clustering, *Microcomputer Information*, 25, 124-125+131, 2009.

536      Rossow, W. B., Leonid, C.G.: Cloud detection using satellite measurements of

537      infrared and visible radiances for ISCCP. *Journal of Climate*, 12, 2341-2369,

538      https://doi.org/10.1175/1520-0442(1993)006<2341:CDUSMO>2.0.CO;2, 1993.

539      Sarkar, V.: Optimized Unrolling of Nested Loops, *International Journal of Parallel*

540      *Programming*, 29, 545-581, https://doi.org/10.1023/A:1012246031671, 2001.

541      Solvsteen, C.: Correlation based cloud-detection and an examination of the split-

542      window method, *Proceedings of SPIE - The International Society for Optical*

543      *Engineering,* 86-97, 1995.

544   Amato, U., Antoniadis, A., Cuomo, V., Cutillo, L., Franzese, M., Murino, L., Serio,

545     C.: Statistical cloud detection from SEVIRI multispectral images, *Remote Sensing*

546     *of Environment*, 112, 750–766, https://doi.org/10.1016/j.rse.2007.06.004, 2008.

547   Wang, Z.: CloudSat Project: CloudSat 2B-CLDCLASS-LIDAR product process

548     description and interface control document, *Jet Propulsion Laboratory*, 2019.

549   Yan, J., Guo, X., Qu, J.: An FY-4A/AGRI cloud detection model based on the naive

550     Bayes algorithm, *Remote Sensing for Natural Resources*, 34, 33-42, 2022.

551   Zhang, W., He, M., Mak, M.W.: Cloud detection using probabilistic neural networks,

552     *Geoscience and Remote Sensing Symposium*, IEEE 2373-2375, 2001.

553   Zhang, Y., William, B. R., Andrew, A. L., Valdar, O., Michael, I. M.: Calculation of

554     radiative fluxes from the surface to the top of atmo- sphere based on ISCCP and

555     other global data sets: Refine- ments of the radiative transfer model and the input

556     data, *Journal of Geophysical Research Atmospheres*, 109, 1-27,

557     https://doi.org/10.1029/2003JD004457, 2004.

558   Zhang, Y., Yang, C., Tao, R.: Multi-temporal Cloud Detection Method for Qin- ghai-

559     Tibet Plateau based with FY-4A Data, *Remote Sensing Technology and Application*,

560     35, 389-398, 2020.

561   Zhen, J., Liu, D., Wang Z.: Analysis of global distribution and seasonal variation

562     characteristics of clouds using CloudSat/CALIPSO satellite data, *Meteorological*

563     *Journal*, 76, 420-433, 2018.