1 **A Machine Learning-Based Cloud Detection and Thermodynamic**

2 **Phase Classification Algorithm using Passive Spectral Observations**

3 **Chenxi Wang[1,2], Steven Platnick[2], Kerry Meyer[2], Zhibo Zhang[3], Yaping Zhou[1,2]**

4

5 [1]Joint Center for Earth Systems Technology, University of Maryland Baltimore County,

6 Baltimore, MD, USA

7 [2]Earth Science Division, NASA Goddard Space Flight Center, Greenbelt, MD, USA.

8 [3]Department of Physics, University of Maryland Baltimore County, Baltimore, MD, USA.

9

10    **Abstract**

11    We trained two Random Forest (RF) machine-learning models for cloud mask and cloud

12    thermodynamic phase detection using spectral observations from VIIRS on Suomi NPP (SNPP).

13    Observations from CALIOP were carefully selected to provide reference labels. The two RF

14    models were trained for all-day and daytime-only conditions using a 4-year collocated

15    VIIRS/CALIOP dataset from 2013 to 2016. Due to the orbit difference, the collocated CALIOP

16    and SNPP VIIRS training samples cover a broad viewing zenith angle range, which is a great

17    benefit to overall model performance. The all-day model uses 3 VIIRS infrared (IR) bands (8.6,

18    11, and 12 $\mu$m) and the daytime model uses 5 Near-IR (NIR) and Shortwave-IR (SWIR) bands

19    (0.86, 1.24, 1.38, 1.64 and 2.25 $\mu$m) together with the 3 IR bands to detect clear, liquid water, and

20    ice cloud pixels. Up to 7 surface types, namely, ocean/water, forest, cropland, grassland, snow/ice,

21    barren/desert, and shrubland, were considered separately to enhance performance for both models.

22    Detection of cloudy pixels and thermodynamic phase with the two RF models were compared

23    against collocated CALIOP products from 2017. It is shown that, with a conservative screening

24    process that excludes the most challenging cloudy pixels for passive remote sensing,  the two RF

25    models have high accuracy rates in comparison with the CALIOP reference for both cloud

26    detection and thermodynamic phase. Other existing SNPP VIIRS and Aqua MODIS cloud mask

27    and phase products are also evaluated, with results showing that the two RF models and the

28    MODIS MYD06 optical property phase product are the top 3 algorithms with respect to lidar

29    observations during the daytime. During the nighttime, the RF all-day model works best for both

30    cloud detection and phase, in particular for pixels over snow/ice surfaces. The present RF models

31    can be extended to other similar passive instruments if training samples can be collected from

32   CALIOP or other lidars. However, the quality of reference labels and potential sampling issues

33   that may impact model performance would need further attention.

**1. Introduction**

35   Detection and classification (DC) of atmospheric constituents using satellite observations is

36   often a critical initial step in many remote sensing algorithms. For example, a prerequisite for cloud

37   optical and microphysical property retrievals is identifying the presence of clouds, i.e., a

38   clear/cloudy classification [*Frey et al.*, 2008; *Heidinger et al.*, 2012]. Additionally, characteristics

39   such as cloud thermodynamic phase are needed as they can strongly impact the

40   scattering/absorption properties of cloud droplets/particles [*Pavolonis et al.*, 2005; *Platnick et al.*,

41   2017]. Similarly, current operational aerosol algorithms can only retrieve aerosol optical depth

42   (AOD) for "non-cloudy" pixels since even slight cloud contamination can result in erroneously

43   high retrieved AOD [*Remer et al.*, 2005]. Therefore, errors in detecting and classifying

44   atmospheric components can significantly impact downstream retrieval products and scientific

45   analyses.

46   There are many examples of hand-tuned DC algorithms designed for satellite instruments. For

47   example, the Moderate Resolution Imaging Spectroradiometer (MODIS) has algorithms

48   developed for cloud masking [*Frey et al.*, 2008; *Ackerman et al.*, 2008], cloud thermodynamic

49   phase [*Baum et al.*, 2012; *Marchant et al.*, 2016], aerosol type [*Levy et al.*, 2013; *Sayer et al.*,

50   2014], and snow coverage over land surfaces [*Hall and Riggs*, 2016]. Decision trees or voting

51   schemes involving multiple thresholds are typically used in these hand-tuned algorithms. The

52   decision tree branches, tests, and thresholds are often determined empirically after a tedious hand

53   tuning/testing process based on the developer's experience and access to validation datasets.

54   Further, the branches and thresholds are often very sensitive to the specific instrument (e.g.,

spectral band pass, calibration, noise characteristics, view/solar geometry sampling). Therefore,

an obvious weakness of these hand-tuned methods is that it is challenging and time consuming to

develop algorithms across multiple instruments and to maintain performance for individual

instruments that may have noticeable calibration drifts. Meanwhile, a well-designed hand-tuned

method may have remarkable performance in a specific region and season yet have significant

biases when applied globally and/or annually [*Cho et al.*, 2009; *Liu et al.*, 2010]. Additional

complexities arise when DC problems become more non-linear across large spatial and temporal

scales, and more variables need to be considered. It is difficult to develop and apply a single or a

few decision trees to complicated non-linear problems that are controlled by dozens or more

variables. As expected, a single decision tree can grow very deep and tend to have a highly

irregular structure in order to consider a large number of features (variables) simultaneously,

leading to a significant overfitting effect (i.e., an over-constrained training that makes predictions

too close to the training dataset but fails to predict future observations reliably). For example,

MODIS provides an all-day cloud phase product based only on infrared (IR) observations

(hereafter referred to as IR-Phase [*Baum et al.*, 2012]). Although it can be expected that the tests

and thresholds should vary with satellite viewing geometry [*Maddux et al.*, 2010], full

consideration of viewing geometries, together with the variations of many other factors such as

surface emission, geolocation, and cloud properties, is very challenging based on manual tuning.

As a consequence, it is found that the liquid water and ice cloud fractions from the IR-Phase

product exhibit noticeable view zenith angle (VZA) dependency (see Figure 12). This is an

undesirable but unavoidable artifact since cloud phase statistics should be independent from

solar/viewing geometry. Such VZA dependencies may strongly affect similar products from

77 geostationary imagers because of the fixed VZA-geolocation mapping. Similar artifacts may also

78 impact aerosol type and retrieval products [*Wu et al.*, 2016].

79     In contrast to hand-tuned methods, Machine Learning (ML) based DC algorithms are designed

80 to autonomously find information (e.g., patterns of spectral, spatial, and/or time series) in one or

81 more given datasets and learn hidden signatures of different objects. An obvious advantage of ML

82 models is that the training process is efficient and highly flexible. Manually defined thresholds or

83 matching conditions to expected spectral patterns are no longer needed. Recently, ML models have

84 been utilized in a wide variety of cloud/aerosol related applications, such as cloud detection

85 [*Thampi et al.*, 2017], cirrus detection and optical property retrievals [*Kox et al.*, 2014; *Strandgren*

86 *et al.*, 2017], surface-level PM2.5 concentration estimation [*Hu et al.*, 2017], and automatic ship-

87 track detections [*Yuan et al.*, 2019]. In this paper, we developed two ML-based DC algorithms for

88 detecting cloud and cloud thermodynamic phase for different local times (i.e., daytime and

89 nighttime) with observations from the Visible Infrared Imaging Radiometer Suite (VIIRS) on

90 Suomi NPP (SNPP). The ML models are trained with collocated observations from SNPP VIIRS

91 and Cloud-Aerosol Lidar with Orthogonal Polarization (CALIOP), with CALIOP data used as the

92 reference. In Section 2, we give a brief discussion of the ML models. Data generated for model

93 training and validation will be introduced in Section 3. Details of the model training and evaluation

94 are shown in Section 4. Section 5 discusses the advantages and potential limitations of the present

95 ML models. Conclusions are given in Section 6.

96 **2. Hand-tuned DC methods and Machine Learning Models**

97 **2.1 Hand-tuned DC methods**

98     All DC algorithms with remote sensing observations are based on the underlying physics of

99 the spectral, spatial, and/or temporal structures of specified objects. In hand-tuned DC algorithms,

100    all the physical rules and structures have to be explicitly defined as various tests and thresholds.

101    For example, the MODIS MOD35/MYD35 cloud mask algorithm uses more than 20 tests with

102    visible/near-infrared (VNIR), shortwave-infrared (SWIR), and infrared (IR) observations [*Frey et*

103    *al.*, 2008] that are carefully designed to consider numerous scenarios, including different surface

104    types (e.g., ocean, land, desert, snow, etc.) and local times (day/night). Similar algorithms are

105    designed for aerosol type and cloud thermodynamic phase classifications. As an example, Figure

106    1 illustrates spectral patterns of 5 typical daytime oceanic scenes (pixel types) observed by SNPP

107    VIIRS. The spectral pattern of each of the 5 scenes, namely, clear sky, liquid water cloud, ice

108    cloud, dust, and smoke, is averaged by using more than 1,000 pixels with the same type. It is clear

109    that the 5 scenes are different in either reflectance ratios between a given VNIR/SWIR band and

110    the 0.86 $\mu$m band, or brightness temperature differences (BTD) between two IR window bands

111    (Figure 1). Consequently, such spectral features are frequently used to differentiate pixel types in

112    DC algorithms. In addition to spectral patterns, simple methods are developed to take into account

113    spatial information. For example, it is found that cloud reflectance usually has larger spatial

114    variability than aerosols [*Martins et al.*, 2002] and clear sky pixels [*Platnick et al.*, 2017].

115    Therefore, spatial variabilities of VNIR and SWIR reflectance bands are used to differentiate

116    clouds from non-cloudy pixels in the current MODIS clear sky restoral (CSR) algorithm [*Platnick*

117    *et al.*, 2017] and Dark Target aerosol retrieval algorithm [*Levy et al.*, 2013].

118    **2.2 Machine learning models**

119        Different from the hand-tuned DC methods, ML algorithms are developed to autonomously

120    learn the hidden spectral/spatial/temporal patterns of different objects. Consequently, manually

121    defined thresholds or matching conditions to expected patterns are no longer needed. In image

122    recognition applications, numerous ML algorithms [e.g., *Joachims* 1998; *Breiman* 1999;

123   *Dietterich* 2000] were developed in late 1990s for independent pixels using a single or small

124   number of decision trees. *Ho* [1998] and many other studies have demonstrated that, although

125   these single or small number of decision trees can always provide maximum prediction accuracies

126   in training processes, significant overfitting effects cannot be avoided. Tremendous efforts have

127   been made to overcome the dilemma between maintenance of prediction accuracy and avoiding

128   overfitting. Among these, the Random Forest (RF) and Gradient Boosting (GB) algorithm

129   [*Breiman* 1999; *Dietterich* 2000; *Friedman* 2001] provide a framework of using a large number of

130   decision trees (ensemble) but a subset of features in each tree to achieve optimization in the

131   performance. It has been demonstrated that the ensemble-based algorithms can largely correct

132   mistakes made by individual trees [*Ji and Ma*, 1997; *Tumer and Ghosh*, 1996; *Latinne et al.*, 2001]

133   and avoid overfitting [*Freund et al.*, 2001]. Currently, the RF and GB algorithms are frequently

134   used in non-linear classification and regression problems. For example, RF models have been used

135   in several cloud/aerosol remote sensing applications, such as differentiating cloudy from clear

136   footprints for the Clouds and the Earth's Radiation Energy System (CERES) instrument [*Thampi*

137   *et al.*, 2017], estimating surface-level PM2.5 concentrations [*Hu et al.*, 2017], and detecting low

138   clouds with the Advanced Baseline Imager (ABI) on the recent Geostationary Operational

139   Environmental Satellites (GOES) [*Haynes et al.*, 2019].  In our study, we also choose the RF model

140   based on its proven record in earth science applications.

141       In the RF model, a final prediction is made based on majority vote computed from probability

142   ($P_i$) of each class ($i^{th}$):

143   $$P_i = \frac{w_i N_i}{\sum_{j=1}^{j=m} w_j N_j},$$   (1)

144 where $m$ is the total number of classes, $N_i$ and $N_j$ are the number of trees that predict the $i$th and $j$th

145 classes, and $w_i$ and $w_j$ are weightings for the $i$th and $j$th classes, respectively. If all trees are equally

146 weighted, $w$ for each individual class is equal to 1. The two most important parameters for tuning

147 the RF algorithm are the number of decision trees ($N_{Tree}$) and the maximum tree depth ($N_{Depth}$).

148 However, an optimal definition of these two parameters is still an open question [*Latinne et al.*,

149 2001]. Larger $N_{Tree}$ and $N_{Depth}$ provides more accurate predictions at the cost of significantly

150 increased computational resources. For many cases, larger $N_{Depth}$ may cause overfitting effects

151 [*Oshiro et al.*, 2012; *Scornet*, 2018]. Generally, the two parameters have to be large enough to let

152 the decision trees have a relatively wide diversity and capture the hidden patterns. However, for

153 practical purposes, the two parameters have to be small enough to prevent the models from

154 overfitting and to reduce computing burden [*Latinne et al.*, 2001; *Scornet* 2018].

155    In this study, we adopt a widely applied RF algorithm in the Scikit-learn Machine Learning

156 package [*Pedregosa et al.*, 2011]. We train two RF models for object DC using SNPP VIIRS

157 spectral observations at two observational times: an all-day RF model using three VIIRS thermal

158 IR observations (hereafter referred to as the RF all-day model) and a daytime-only RF model that

159 uses both VNIR/SWIR and thermal IR observations (hereafter the RF daytime model). The models

160 are trained to detect clear sky, liquid water cloud, and ice cloud pixels with single pixel level

161 information. Parameters of the two RF models will be tuned and tested carefully to achieve the

162 best accuracy and to avoid the overfitting effect. Details will be discussed in Section 4.

163 **3. Data**

164 **3.1 Reference label of pixels**

165    Space-borne active sensors, such as CALIOP onboard CALIPSO [*Winker et al.*, 2013], the

166 Cloud-Aerosol Transport System (CATS) [*McGill et al.*, 2015] onboard the International Space

167 Station (ISS), and CPR on board CloudSat [*Stephens et al.*, 2002], are frequently used to evaluate

168 the performance of hand-tuned cloud/aerosol DC and property retrieval algorithms designed for

169 passive sensors [*Stubenrauch et al.*, 2013; *Wang et al.*, 2019]. CALIPSO, a key member of the

170 Afternoon Constellation of satellites (A-Train) until its exit on 13 September 2018 to join CloudSat

171 in a lower orbit, began providing profiling observations of the atmosphere in 2006 [*Winker et al.*,

172 2013]. The CALIPSO lidar CALIOP operates at wavelengths of 532 nm and 1064 nm, measuring

173 backscattering profiles at a 30-meter vertical and 333 m along-track resolution. CALIOP also

174 measures the perpendicular and parallel signals at 532 nm, along with the depolarization ratio at

175 532 nm that is frequently used in cloud phase discrimination algorithms because of its strong

176 particle shape dependence. The CALIOP Version 4 Level 2 1 km/5km Layer product is used to

177 provide reference cloud phase labels in both model training and validation stages.

178 While the CATS lidar and the CloudSat radar CPR also provide profiling information, both

179 have limitations that preclude their use here. CATS had a relatively short life time (from January

180 2015 to October 2017), and its low inclination angle (51°) orbit aboard the ISS excludes sampling

181 of high-latitude regions [*Noel et al.*, 2018]. CloudSat CPR observes reflectivity profiles at 94-GHz,

182 which are more sensitive to optically thicker clouds consisting of large particles but are blind to

183 aerosols and optically thin clouds. CloudSat also has difficulty in detecting clouds near the surface

184 due to the surface clutter effect [*Tanelli et al.*, 2008]. Therefore, only CALIOP data are used to

185 provide reference cloud phase labels in this study.

186 **3.2 RF model input**

187 It should be pointed out that ML models use similar input datasets as hand-tuned methods. The

188 input variables (features) and reference labels of the present RF models are carefully selected based

189 on prior physical knowledge of the spectral characteristics of each object.

9

190    VIIRS on SNPP and the NOAA-20+ series provides spectral observations from 0.4 to 12 $\mu$m

191    at sub-kilometer spatial resolutions [*Lee et al.*, 2006]. Specifically, VIIRS has 16 moderate

192    resolution bands (M band) and 5 higher resolution imagery bands (I band) at 750 m and 375 m

193    nadir resolutions, respectively. The spectral capabilities of VIIRS allow for extracting abundant

194    information on the surface and atmospheric components, such as clouds [*Ackerman et al.*, 2019]

195    and aerosols [*Sayer et al.*, 2017]. It is also worth noting that VIIRS utilizes an on-board detector

196    aggregation scheme that minimizes pixel size growth in the across-track direction towards swath

197    edge [*Cao et al.*, 2013]. As an example, although the VIIRS M-bands and MODIS 1 km bands

198    have similar nadir spatial resolutions, the VIIRS across-track pixel size increases to roughly

199    1.625 km at scan edge, which is much smaller than a MODIS pixel size of roughly 4.9 km at scan

200    edge [*Justice et al*., 2011]. Another obvious advantage of using SNPP VIIRS rather than Aqua

201    MODIS data is that, due to the CALIPSO and SNPP orbit differences, the training samples cover

202    a broader viewing zenith angle range, which is a great benefit to overall model performance.

203    Consequently, Level-1B M-band observations from the SNPP VIIRS are used here.

204    Ancillary data, including the surface skin temperature, spectral surface emissivity, surface

205    types, and snow/ice coverage, are important in cloud DC related remote sensing applications [*Frey*

206    *et al.*, 2008; *Wolters et al.*, 2008; *Baum et al.*, 2012] and cloud/aerosol retrievals [*Levy et al.*, 2013;

207    *Wang et al.*, 2014; 2016a; 2016b; *Meyer et al.*, 2016; *Platnick et al.*, 2017]. The inst1_2d_asm_Nx

208    product (version 5.12.4) from the Modern-Era Retrospective Analysis for Research and

209    Applications, Version 2 (MERRA-2) [*Gelaro et al.*, 2017] is utilized to provide the hourly

210    instantaneous surface skin temperature and 10-meter surface wind speed. The UW-Madison

211    baseline fit land surface emissivity database [*Seemann et al.*, 2008] and the Terra/Aqua MODIS

212    combined Land surface product (MCD12C1 [*Sulla-Menashe and Friedl* 2018]) are used to provide

213    monthly mean land surface emissivities for the mid-wave to thermal IR bands (3.6 ~ 14.3 $\mu$m) and

214    surface white sky albedo for the VNIR bands (0.4 ~ 2.3 $\mu$m), respectively, at a 0.05×0.05° spatial

215    resolution. Surface types and snow/sea ice coverage data are from the International Geosphere-

216    Biosphere Programme (IGBP) and daily Near-real-time Ice and Snow Extent (NISE) data [*Brodzik*

217    *and Stewart*, 2016], respectively.

**3.3 Clear and cloud phase classifications from existing VIIRS and MODIS products**

219    Since the present RF models are trained with SNPP VIIRS observations, the first priority of

220    this study is evaluating and comparing the trained RF models with CALIOP and the existing VIIRS

221    cloud products. However, existing cloud mask and phase products from Aqua MODIS are still

222    used as a reference in this work.

223    The Aqua MODIS and SNPP VIIRS CLDMSK (cloud mask) and CLDPROP (cloud top and

224    optical properties) [*Ackerman et al.*, 2019] products represent NASA's effort to establish a long-

225    term consistent cloud climate data record, including cloud detection and thermodynamic phase,

226    across the MODIS and VIIRS observational records. While the CLDMSK (version 1.0) and

227    CLDPROP (version 1.1) algorithms share heritage with the standard Collection 6.1 MODIS cloud

228    mask (MYD35) and cloud top and optical properties (MYD06) algorithms, the algorithms use only

229    a subset of bands common to both sensors to minimize differences in instrument spectral

230    information content.

231    The CLDMSK and MYD35 algorithms use a variety of band combinations and thresholds

232    depending on cloud and surface types [*Frey et al.*, 2008; *Ackerman et al.*, 2008]. Meanwhile, the

233    algorithms use different approaches for daytime (i.e., solar zenith angle less than 85°) and

234    nighttime pixels. In the CLDMSK and MYD35 algorithms, pixels are categorized into four

235  categories, namely confident clear, probably clear, probably cloudy, and cloudy. The CLDPROP

236  and MYD06 algorithms separate cloudy and probably cloudy pixels into liquid water, ice, and

237  unknown phase categories. Specifically, the MYD06 product includes two cloud phase algorithms:

238  an IR-Phase algorithm [*Baum et al.*, 2012] that uses observations in four MODIS IR bands for

239  daytime and nighttime phase classification (hereafter referred to as the MYD06 IR-Phase), and a

240  daytime-only algorithm designed for the cloud optical properties retrievals [*Marchant et al.*, 2016;

241  *Platnick et al.*, 2017] that uses VNIR/SWIR and IR observations (hereafter referred to as the

242  MYD06 OP-Phase). A notable change for the VIIRS/MODIS CLDPROP algorithm with respect

243  to the standard MODIS MYD06 algorithm is the replacement of the MYD06 IR-Phase by a NOAA

244  operational algorithm originally developed for Clouds from AVHRR-Extended (CLAVR-x)

245  [*Heidinger et al.*, 2012] and now applied to VIIRS. This algorithm is used to provide cloud top

246  properties, including thermodynamic phase (hereafter CLDPROP CT-Phase), in the absence of the

247  MODIS $CO_2$ IR gas absorption bands. IR bands are primarily used in the CLDPROP CT-Phase

248  algorithm, while complementary SWIR bands are used when available. The MYD06 OP-Phase

249  algorithm, applied to daytime pixels only, is included with only minor alteration (related to cloud

250  top properties changes) in the VIIRS/MODIS CLDPROP product (hereafter referred to as the

251  CLDPROP OP-Phase).

252  Although the MYD06 and CLDPROP OP-Phase products are developed for "cloudy" and

253  "probably cloudy" pixels from the MYD35 and CLDMSK products, a Clear Sky Restoral (CSR)

254  algorithm [*Platnick et al.*, 2017] is implemented to remove "false cloudy" pixels from the clear-

255  sky conservative MYD35 and CLDMSK products. Specifically, the CSR uses a set of spectral and

256  spatial reflectance variability tests to remove dust, smoke, and strong sunglint pixels that are

257  erroneously identified as "cloudy" or "probably cloudy" by the MYD35 and CLDMSK products

258　[*Platnick et al.*, 2017]. One should keep in mind that the CSR algorithm is only applied for the

259　optical property retrievals. Thus, the MYD35 and CLDMSK, and consequently the MYD06 IR-

260　Phase and CLDPROP CT-Phase, may have "false cloudy" pixels in comparison with CALIOP,

261　while the impact on the MYD06 and CLDPROP OP-Phase is reduced due to the CSR algorithm.

262　The cloud mask and thermodynamic phase products used in this study are summarized in Table 1.

263　**4. Model training and validation**

264　　Here we discuss the training of the all-day and daytime RF models for different surface types.

265　Both shortwave (SW) and IR observations will be used in the daytime models while only IR

266　observations will be used in the all-day models. ML model performance is strongly dependent on

267　the quality of training samples. In this study, the two RF models are trained and tested with simple

268　yet highly confident samples (Section 4.2). With this training strategy, the RF models are expected

269　to capture the key spectral features from the pure samples efficiently. As discussed in Section 4.4,

270　we conducted a model validation that evaluates performance of the two models for simple cases.

271　Furthermore, an analysis of probability distributions from the RF all-day model is conducted to

272　demonstrate that the RF models have capability to recognize spectral features from more than one

273　category when atmospheric columns are more complicated.

274　**4.1 Surface Types**

275　　RF models are trained for different surface types, defined here by the Collection 6 (C6) MODIS

276　annual IGBP surface type product (MCD12C1), to improve model performance over a single

277　general model for all surface types. Although the MCD12C1 product includes up to 18 surface

278　types, for this work we attempt to reduce the total number of surface types by combining surface

279　types with similar spectral white sky albedos and emissivities, as suggested by *Thampi et al.*

280　[2017]. An annual global IGBP surface type map and surface albedo data from the MODIS

281  MCD12C1 [*Sulla-Menashe and Friedl* 2018] and a UW-Madison monthly global land surface

282  emissivity database [*Seemann et al.*, 2008] are used to generate the climatology of land surface

283  white-sky albedo and IR emissivity spectra. The UW-Madison database is derived using input

284  from the MODIS operational land surface emissivity product MOD11 [*Wan et al.,* 2004] at six

285  wavelengths located at 3.8, 3.9, 4.0, 8.6, 11, and 12 $\mu$m.  A baseline fit method is applied to fill

286  the spectral gaps and provides a more comprehensive IR emissivity dataset at 10 wavelengths from

287  3.6 to 14.3 micron for global land surface with a 0.05° spatial resolution [*Seemann et al.*, 2008].

288  The MODIS MCD12C1 product also provides a white-sky albedo dataset at 0.47, 0.56, 0.66, 0.86,

289  1.24, 1.64, and 2.13 $\mu$m with a 0.05° spatial resolution [*Sulla-Menashe and Friedl* 2018]. The

290  means and standard deviations of surface emissivity and white-sky albedo spectra are shown in

291  Figures 2 a) and 3 a), respectively, for 16 different land surface types generated from the UW-

292  Madison and MCD12C1 data in 2015. Land surface types with similar IR emissivity and SW

293  white-sky albedo spectra are grouped to reduce to the total number of land surface types to 6

294  (forest, cropland, grassland, snow/ice, barren/desert, and shrubland), as shown in Figures 2 (b-f)

295  and 3 (b-f). Figure 4 shows an example map of the reduced global surface type data generated

296  from the MCD12C1 product for 2015.

297  **4.2 Generating Training/Validation Datasets**

298     The training and validation data are obtained from a 5-year (2013-2017) SNPP VIIRS and

299  CALIOP collocated dataset. The collected dataset is generated with a collocation algorithm that

300  fully considers the spatial differences between the two instruments and parallax effects, as

301  described in *Holz et al.* [2008]. The SNPP VIIRS data include L1B calibrated reflectance and

302  brightness temperatures, and the CALIOP data include the L2 1km/5km cloud and aerosol layer

303  products. Although more than 332 million VIIRS 750m pixels are collocated with CALIOP

304    observations, 130.6 million of these pixels (39.3%) that include only aerosol-free, homogeneous,

305    clear (39.1 million) or single-phase cloud (49.7 million liquid and 41.8 million ice) pixels are used

306    in our training/validation process. Unless otherwise specified, "*aerosol-free*" is defined as those

307    pixels having collocated CALIOP 5km column 532 nm aerosol optical depth less than 0.05,

308    "*homogeneous*" is defined as those pixels for which the collocated CALIOP 1km and 5km

309    products have the same pixel labels, and "*single-phase cloud*" is defined as those pixels for which

310    the collocated CALIOP 1km and 5km products indicate the same thermodynamic phase for all

311    identified cloud layers. More details are given in Table 2.

312        A strict three-step quality control process is applied to collect samples for the

313    training/validation process. First, VIIRS 750 m pixels that are potentially contaminated by aerosol

314    are excluded using a threshold of 0.05 column AOD at 532 nm from the CALIOP L2 5 km aerosol

315    layer product. Second, each aerosol-free pixel is labelled by one of four categories, namely, "clear

316    sky" and "liquid-water cloud", "ice cloud", and "ambiguous" with the CALIOP L2 1km/5km layer

317    product. The "ambiguous" pixels, including uncertain/unknown cloud phases from CALIOP

318    and/or overlapping objects belonging to different types (e.g., cirrus over liquid), are discarded.

319    Third, horizontally inhomogeneous pixels, determined when the CALIOP 1km label changes

320    within 5 consecutive VIIRS pixels, or pixels with inconsistent CALIOP 1km and 5km labels, are

321    discarded. Figure 5 shows the global distributions of the 5-year collocated clear (first row) and

322    cloudy pixels (second row) before and after applying the three-step quality control. Globally, 50%

323    of all clear pixels are excluded due to contamination of broken-cloud and/or aerosol. In particular,

324    a large fraction of clear pixels in central Africa, India, and southern China (Figure 5c) are excluded

325    due to relatively large aerosol optical thicknesses in those regions. About 40% of global cloudy

326    pixels (Figure 5f) are excluded due to cloud heterogeneity and aerosol contamination. The

327　minimum selection rate (~20%) can be found in some particular regions, such as the Inter Tropical

328　Convergence Zone (ITCZ), where clouds have complicated horizontal/vertical structures due to

329　strong convections (i.e., clouds are highly heterogeneous in both the horizontal and vertical

330　dimensions). The remaining data are separated into a training/testing population that consists of

331　32.4, 41.2 and 34.9 million pixels for clear sky, liquid water cloud, and ice cloud from years 2013-

332　2016, respectively, and a validation dataset that consists of 6.9, 8.5 and 7.0 million pixels of clear-

333　sky, liquid water cloud, ice cloud, respectively from year 2017.

334　**4.3 RF model training and configuration**

335　　RF model performance is determined by both its inputs (spectral or other information) and its

336　configuration ($N_{Tree}$ and $N_{Depth}$). Therefore, extensive testing must be conducted to find the optimal

337　inputs and configuration. The 4-year collocated VIIRS-CALIOP dataset from 2013 to 2016 after

338　quality control (see Section 4.2) is used for both training (75%) and testing (25%) purposes. The

339　testing set, also known as cross-validation set, is used to tune and optimize the RF model

340　parameters. Here we define an accuracy score to evaluate the overall model performance. The

341　accuracy score is the ratio of pixels (samples) where both the CALIOP and RF model have the

342　same categories to total pixels. In this study, we tested six groups of input variables for each RF

343　model. The set of model input variables with a relatively high accuracy score and low

344　memory/computing requirement will be selected.

345　　Table 3 provides accuracy scores of the IR-based all-day model trained and tested with

346　different inputs. It shows that with a fixed RF model configuration ($N_{Tree}$ = 150 and $N_{Depth}$ = 15),

347　the RF all-day model with input #4 and #6 have the best overall accuracy scores for all surface

348　types. Generally, by including surface skin temperature ($T_s$) and geolocation (i.e., latitude and

349　longitude), the accuracy scores for all surface types increase by 2-3%. The surface emissivity

350 vector $\varepsilon_s$ is less important, likely because this information is highly correlated to surface type and

351 geolocation. In this study, input #4 is selected mainly because with similar performance, it requires

352 less memory and computing resources, and it is quite possible that more uncertainty is introduced

353 with the use of a surface emissivity vector $\varepsilon_s$ from another retrieval product.

354     A set of model configurations ($N_{Tree}$ and $N_{Depth}$) are also tested based on the selected input #4.

355 While the number of trees and the maximum depth of individual trees are important determinants

356 for RF model performance, the overall accuracy scores for all surface types are less sensitive to

357 these two model parameters when more than 100 trees and 10 maximum tree depths are used (not

358 shown here). Therefore, we trained the RF all-day models with input #4 and the model

359 configuration used in Table 3, i.e., $N_{Tree}$ = 150 and $N_{Depth}$ = 15.

360     Similar input variable tests for the RF daytime model (IR plus NIR and SWIR observations)

361 showed that the optimal input includes reflectances in the 0.86, 1.24, 1.38, 1.64 and 2.25$\mu$m bands,

362 BTs in the same 3 IR bands used in the all-day model, geolocation, and solar/satellite viewing

363 zenith angles (See Table 4). The same model configuration used in the all-day model, e.g., 150

364 trees with the maximum depth 15, is used in the daytime model. The accuracy scores of the RF

365 daytime model are higher than the RF all-day model by 2-3% over almost all surface types except

366 high-latitude regions covered by snow and ice, where the daytime model accuracy score is higher

367 by up to 6% than the all-day model due to the inclusion of the 1.38, 1.64 and 2.25µm SWIR bands.

368 **4.4 Evaluating the RF Models**

369     The trained RF all-day and daytime models are validated using collocated CALIOP data in

370 2017. Existing VIIRS cloud products CLDMSK and CLDPROP (see Table 1) are included for

371 direct comparison with the RF models and CALIOP reference. Several other products, such as the

372    MODIS CLDMSK and CLDPROP and standard MYD35 and MYD06, are also included for

373    comparison although they could be different from the RF models due to other non-algorithm

374    reasons, such as the VZA and pixel size differences mentioned before.

375    *4.5.1 Cloud mask*

376    Cloud mask from the two RF models and VIIRS/MODIS products are first compared with

377    CALIOP lidar observations. For the two models, a cloudy pixel indicates a predicted label "liquid"

378    or "ice". Here we define cloudy and clear pixels as "positive" and "negative" events, respectively.

379    A true positive rate (TPR) and false positive rate (FPR) can then be used to evaluate model

380    performance. The TPR and FPR are defined as:

$$\text{TPR} = \frac{TP}{TP+FN}, \tag{2}$$

381

$$\text{FPR} = \frac{FP}{FP+TN}, \tag{3}$$

382

383    where TP (True Positive) and TN (True Negative) are the number of lidar-labeled "cloudy" and

384    "clear" pixels, respectively, that are correctly detected by the models; whereas FN (False Negative)

385    and FP (False Positive) are the number of lidar-labeled "cloudy" and "clear" pixels incorrectly

386    identified by the models. Therefore, TPR, also called model sensitivity, indicates the fraction of

387    all positive events (i.e., lidar cloudy pixels) that are correctly detected by the models. Similarly,

388    FPR, also called false alarm rate, indicates the fraction of all negative events (i.e., lidar clear pixels)

389    that are incorrectly detected as positive (cloudy). TPR and FPR are two critical parameters in

390    model evaluation. A perfect model is associated with a high TPR (close to 1) and a low FPR (close

391    to 0).

392    Figure 6 shows daytime cloud mask TPR-FPR plots from the two RF models and the other

393    products listed in Table 1. Globally, all products agree well with lidar observations (Figure 6a).

18

394    The overall TPRs are higher than 0.94 and FPRs are lower than 0.08. The RF daytime model (red

395    circle), with a TPR of 0.97 and an FPR of 0.05, is slightly better than the RF all-day model (yellow

396    circle) and other products. Figure 6b-6h show comparisons over different surface types. It is clear

397    that the RF daytime model has a robust performance for all surface types. The MODIS MYD35

398    cloud mask algorithm (black circle) performs best over ocean but has a relatively high FPR (0.22)

399    over forest and low TPR over snow/ice and barren (0.85) regions. As mentioned in Section 3, the

400    "false cloudy" pixels from MYD35 and CLDMSK may increase the FPRs correspondingly.

401    The RF all-day model works fairly well and is comparable to other products for all surface

402    types regardless of the fact that it only uses three IR window channels from VIIRS while all other

403    products in the daytime models use VNIR observations. Nighttime (SZA > 85°) cloud mask

404    comparisons are shown in Figure 7. The overall performances of all operational products decrease

405    in particular for snow/ice regions. For example, the VIIRS/MODIS CLDMSK products over

406    snow/ice surface have large fractions of missing "cloudy" pixels (e.g., TPRs < 0.7) and false alarm

407    rates (FPRs > 0.2) over snow/ice surface. The decrease is more likely explained by the lack of

408    SWIR bands and the small cloud-snow/ice surface temperature contrast during the nighttime of

409    summer polar regions. However, the RF all-day model has the best performance for nighttime

410    pixels, indicating the strong capability of ML based algorithm in capturing hidden spectral features

411    and optimizing dynamic thresholds of clear and cloudy pixels.

412    *4.5.2 Cloud thermodynamic phase*

413    The RF cloud thermodynamic phase products are also compared with CALIOP lidar and

414    existing VIIRS and MODIS products. For consistent nomenclature, we arbitrarily define ice clouds

415    and liquid water clouds as "positive" and "negative" events, respectively. A low TPR indicates

416    underestimation of ice cloud fraction, while a high FPR indicates a large fraction of liquid water

417  cloud samples are identified as ice cloud. To focus on cloud thermodynamic phase classification,

418  pixels detected as "clear" by either the lidar reference labels or by the RF models and existing

419  products are excluded. The OP-Phase from both MYD06 and CLDPROP, and the IR-Phase from

420  MYD06, have an "unknown phase" category, which is not included in the TPR-FPR analysis.

421  Figure 8 shows daytime cloud phase TPR-FPR plots from the two RF models and the

422  MODIS/VIIIRS products. The two RF models and the MODIS MYD06 OP-Phase are the top 3

423  phase algorithms for all surface types. The MODIS MYD06 IR-Phase, MODIS/VIIRS CLDPROP

424  OP-Phase, and CT-Phase have either relatively lower TPRs or higher FPRs over particular surface

425  types, such as shrubland, snow/ice, and barren regions. Comparisons between nighttime phase

426  algorithms are shown in Figure 9. For nighttime clouds, the RF all-day model works better than

427  both CT-Phase and IR-Phase algorithms for all surface types. Overall, the performance of the

428  hand-tuned algorithms decreases significantly over snow/ice or barren surfaces. For example, the

429  TPR-FPR plot shows that over daytime snow/ice surface (Figure 8 g), the MODIS CLDPROP OP-

430  Phase and MODIS MYD06 IR-Phase frequently predict liquid water cloud as ice cloud. Similar to

431  the daytime plot, the MYD06 IR-Phase also shows a high FPR rate over snow/ice surface,

432  indicating an overestimated (underestimated) ice (liquid water) cloud fraction. Possible reasons

433  include strong surface reflection, low surface cloud contrast, relatively less training samples and

434  high solar zenith angles. However, the two RF models work fairly well and show consistent

435  accuracy rates across all surface types.

436  It is also important to note that the number of pixels used for cloud phase TPR-FPR

437  comparisons in Figures 8 and 9 are different for products that have "unknown phase" categories,

438  namely, MYD06 IR-Phase, MYD06 OP-Phase, and CLDPROP OP-Phase. As shown in Table 5,

439  the MYD06 IR-Phase has a relatively large "unknown phase" phase fraction (15% for all surface

440    types and 34% for snow/ice) in comparison to the OP-Phase products from both MYD06 and

441    CLDPROP, which have 2~3% "unknown phase" fraction approximately.

442    As discussed in Section 2.2, recall that the RF model predicted pixel type is derived by setting

443    thresholds on the probabilities for each classification type, e.g., an ice phase decision is reached if

444    the probability of ice is greater than the probabilities of liquid and clear. Figure 10 shows the

445    probability distribution functions of the RF all-day model for four scene types as determined by

446    collocated CALIOP, namely, (a) clear, (b) liquid, (c) ice, and (d) multi-layer clouds with different

447    thermodynamic phases (e.g., ice over liquid). As expected, for the first three types, which are

448    included in the training/validation processes, the probability distributions have strong peaks close

449    to either 0 or 1. For the multiple phase cases (panel d), the liquid and ice probabilities are more

450    broadly distributed, indicating that the model may recognize signals from both liquid and ice and

451    therefore provide ambiguous phase results. More nuanced thresholds can therefore be applied to

452    the probabilities, for instance to create an "unknown" phase category following MYD06 and

453    CLDPROP convention [*Marchant et al.*, 2016] that can indicate complicated cloud scenes.

454    Furthermore, the probabilities themselves can provide a useful quality assurance metric for

455    downstream cloud property retrievals that often must make an assumption on cloud phase.

456    Nevertheless, assigning an appropriate phase for downstream imager-based cloud property

457    retrievals is difficult for complex, multilayer cloud scenes, as such an assignment often depends

458    on the optical/microphysical properties and vertical distribution of the cloud layers in the scene

459    [*Marchant et al.*, 2020]. Further investigation is necessary to understand how to use the RF phase

460    probabilities more quantitatively in complicated cases.

461    Figure 11 shows monthly mean daytime cloud and phase fractions from the VIIRS CLDMSK

462    and CLDPROP OP-Phase products (top row), and those from the RF daytime model (second row),

463    in January 2017. For the cloud mask comparison, cloud fractions (CF) from the two products have

464    similar spatial patterns, while it is also clear that the VIIRS CLDMSK CFs are higher over tropical

465    oceans by approximately 10% and lower over land by 5% (Figure 11 c). This is consistent with

466    the cloud mask TPR-FPR analysis shown in Figure 6. Over the tropical ocean, the VIIRS

467    CLDMSK is more "cloudy", probably due to a fraction of sunglint pixels that are detected as liquid

468    clouds, leading to a large FPR rate. Another reason for the relatively large cloud fraction (or liquid

469    water cloud fraction) difference is that in regions covered by "broken" cumulus clouds, and or

470    clouds with more complicated structures, the inherent viewing geometry differences in the training

471    datasets may adversely affect the performance of the RF models. For example, CALIOP, with a

472    nadir viewing geometry may observe clear gaps between two small cloud pieces, while VIIRS,

473    with an oblique viewing angle, detects broken liquid clouds nearby or high clouds along its long

474    line-of sight. Comparison between the VIIRS product and the RF daytime model shows more ice

475    clouds from the RF daytime models over land, which is consistent with the cloud phase TPR-FPR

476    plots as shown in Figure 8. The RF daytime model may have better performance due to the

477    consideration of surface type. However, it is also important to notice that due to the lack of

478    "aerosol" types in current training, in central Africa, the RF models may misidentify elevated

479    smoke as ice cloudy pixels. For most land surface types except snow/ice, the CLDPROP OP-Phase

480    has lower TPR rates than the RF daytime models by 0.1, in comparison with the CALIOP.

481       In addition to the higher CFs over low latitude ocean from the VIIRS CLDMSK product, more

482    pronounced CF (liquid) differences can be found in northeast and northwest China. Cloud

483    differences in the two regions are spatially correlated with locations that have heavy aerosol

484    loadings or snow coverage. For example, heavy aerosol loadings due to pollution in Northeast

485    China, and a wide land snow coverage in Northwest China are frequently observed in the winter.

486    The VIIRS CLDMSK may identify pixels with white surface and heavy aerosol loadings as

487    "cloudy". Some of these pixels are expected to be restored to clear-sky category in the CLDPROP

488    OP-Phase product (Figure 11 f and i). As evidence, Figure 12 shows comparisons between the

489    VIIRS products and the RF daytime model in July 2017. The large cloud (liquid) fraction

490    differences over North China vanish in the summer. This indicates that the RF models might be

491    able to handle complicated (or unexpected) surface type and strong aerosol events better than the

492    hand-tuned VIIRS algorithm. However, further investigation is required to understand the

493    performances of both the VIIRS products and the RF models.

494    **5. Discussion**

495        In this Section, we will review the strengths and potential limitations and weaknesses of the

496    RF models.

497    **5.1 Advantages**

498        The above results show that, for the screened clear/cloudy samples, the two RF models have

499    better and more consistent performance over different regions and surface types in comparison

500    with the MODIS and VIIRS products, suggesting the potential to improve the overall performance

501    in more global operational applications. In addition to better performance, it is convenient and

502    efficient to apply the present RF models or other similar ML-based models to other instruments

503    similar to VIIRS, such as the geostationary imagers Advanced Himawari Imager (AHI) on

504    Himawari-8/9, the ABI on GOES-16/17, and the Spinning Enhanced Visible and Infrared Imager

505    (SEVIRI) on Meteosat Second Generation, as long as reliable reference pixel labels are available.

506    With hand-tuned methods, adjustment is always required in the case of calibration changes,

507    algorithm porting to another similar instrument, or changes in solar/viewing geometries and

508    surface conditions. Manual adjustments can be time-consuming (e.g., months or years), whereas

509    the two RF models used in this study were trained and tested for 7 surface types and using different

510    input variables in 3 hours (on an HPC Platform using 32 Intel Xeon Gold 6126 Processors @ 2.60

511    GHz).  More important, manual algorithm adjustment may not provide the best continuity between

512    two instruments. For example, although the MODIS CLDPROP OP-Phase and VIIRS CLDPROP

513    OP-Phase are designed for climate record continuity purpose, cloud thermodynamic phases from

514    the two products are different by up to 4% for all surface pixels, and by up to 10% over surfaces

515    covered by snow/ice (see Figure 8 light blue and light green dots). Further investigation is

516    necessary to understand if, using ML approaches, a better climate record continuity will be

517    achieved with a uniform training dataset. Besides providing a discrete category for each pixel, the

518    RF models provide an ensemble of predictions and probabilities of individual categories, which

519    are useful diagnostic variables in evaluating models in complicated scenarios.

520    **5.2 Limitations and possible caveats**

521        Although the evaluation demonstrates that the current RF models are highly consistent with

522    CALIOP, the models may suffer some artifacts due to the quality of the training data and due to

523    sampling issues.

524    *5.2.1 Quality of the training/validation data*

525        The RF models learn spectral structures of cloud/clear pixels according to the reference labels.

526    As a consequence, the present model performance relies heavily on the quality of CALIOP Level-

527    2 data. It is already known that the lidar signal has limitations in detecting the bottom of an

528    optically thick cloud or lower level clouds underneath an opaque cloud [*Sassen and Cho*, 1992].

529    Some complicated multiple-phase scenes may be misidentified as simple single-phase scenes due

530    to the penetration limit of CALIOP (e.g., the uppermost ice cloud optical thickness greater than 3).

531    Using combined CALIOP and CloudSat data as reference in the future could be a better way to

532    improve the training/validation datasets [*Marchant et al.*, 2020]. However, as noted in that study,

533    CloudSat observations cannot be used without careful filtering since a multilayer scene that is

534    radiatively indistinct from the upper level cloud layer is not necessarily consistent with multilayer

535    detection detected from a cloud radar.

536    Additional uncertainties may come from the inconsistency in view angles between the

537    collocated CALIOP labels and VIIRS spectral observations. For instance, CALIOP always has a

538    quasi-nadir viewing angle (e.g., 3°) whereas the collocated VIIRS observations have a wide VZA

539    range (e.g., 0° to 50°). A wide VIIRS VZA range in the training dataset improves model

540    performance, especially for predicting VIIRS pixels with large VZAs. However, the difference

541    between the CALIOP and VIIRS viewing geometry could create undesirable artifacts in the

542    training process. As shown in Figure 11, in the descending areas of the Hadley cell over low-

543    latitude ocean, where marine boundary layer clouds are dominant, there are relatively large CF

544    differences between the CLDMSK and the RF models. A reason for the large liquid cloud fraction

545    differences is that the quality of training datasets decreases in regions covered by "broken"

546    cumulus clouds, and or clouds with more complicated structures. Further investigation is required

547    to check if the training dataset collection process introduces sampling bias into the training dataset.

548    *5.2.2 Sampling issue*

549    Uneven sampling may also influence the training of RF models. Figure 13 shows the cloud

550    fraction as a function of viewing geometry. Quasi-constant fractions of both liquid and ice clouds

551    are found for all operational products and the RF models when VZAs are smaller than 45°, except

552    the MODIS MYD06 IR-Phase, which has a strong VZA dependency. However, liquid (ice) cloud

553    fractions from the two RF models increase (decrease) rapidly at high VZAs (greater than 50°),

554    which is likely caused by the sampling issue. A significant fraction of the training data (greater

555    than 98%) is located in the region with VZA less than 50° (see the gray dashed distributions in

556    Figure 13). It is difficult to mitigate this issue using collocated VIIRS-CALIOP data or

557    observations from other similar instruments in the training process. One possible way is using

558    model-generated synthetic training data and labels with reliable radiative transfer models. Results

559    from the RF daytime model are not shown in Figure 13 since they are highly consistent with the

560    RF all-day model.

561    *5.2.3 Labeling strategy*

562    For RF or other ML models, each pixel's classification is determined by prediction

563    probabilities ($P$) of all potential types. Here we selected a regular strategy that labels a pixel using

564    the class with the highest probability (see Eq. 1). This strategy is logical for problems with two

565    categories (e.g., cloud mask only). For problems including 3 or more classes, however, the present

566    strategy is not the only way to label pixels. For example, a pixel is labeled as "clear" if $P_{clear}$ is

567    larger than both $P_{liquid}$ and $P_{ice}$ according to the current labeling strategy. It is also possible that,

568    for the same pixel (less than 0.5% for the two RF models), $P_{clear}$ is lower than the sum of $P_{liquid}$

569    and $P_{ice}$, making a "cloudy" label more appropriate. For the cloud mask and phase problem

570    discussed in this paper, in addition to pixel labels, users must be aware of probabilities of the three

571    types. Another possible way to avoid the ambiguous labeling is using two RF models, one for

572    cloud masking and one for phase, such that a "clear" or "cloudy" label is given first by the cloud

573    mask model, while a corresponding "liquid" or "ice" label is assigned to "cloudy" pixels in the

574    cloud phase model. However, two RF models double the training process and require more

575    computing resources in operational applications.

576    **6. Conclusions**

577    Two Machine-Learning Random Forest (RF) models were trained to provide pixel types (i.e.,

578    clear, liquid water cloud, and ice cloud) using VIIRS 750-meter spectral observations. A daytime

579    model that uses NIR, SWIR, and IR bands and an all-day model that only uses IR bands were

580    trained separately. In the training processes, reference pixel labels are from collocated CALIOP

581    Level 2 1 km cloud layer and 5 km aerosol layer products from 2013 to 2016. Careful tests were

582    conducted to optimize model input and configuration. The two RF models were trained for 7

583    different surface types (i.e., ocean/water, forest, cropland, grassland, snow/ice, barren/desert, and

584    shrubland) to improve model performance. In addition to geolocation and solar/satellite geometry

585    information, we found that using 5 NIR and SWIR bands (0.86, 1.24, 1.38, 1.64 and 2.25 $\mu$m) and

586    three IR bands (8.6, 11, and 12$\mu$m) in the daytime RF model and using the three IR bands and

587    surface temperatures in the all-day RF model achieved great performances for all surface types.

588    The cloud mask and thermodynamic phase classifications from the two RF models were

589    validated using the selected aerosol-free, homogeneous samples in 2017. For daytime cloud mask

590    comparisons over all surface types, the RF daytime model, with a high TPR (0.93 and higher) and

591    low FPR (0.07 and lower), performs best among all models evaluated, including MODIS MYD35

592    and MODIS/VIIRS CLDMSK products. The RF all-day model works fairly well and is

593    comparable to other products for all surface types, even in daytime when all other products use

594    shortwave observations and it does not. For the nighttime cloud mask, the RF all-day model has

595    the best performance over all products, demonstrating the strong capability of ML-based

596    algorithms for capturing hidden spectral features of clear and cloudy pixels. All nighttime products

597    perform slightly weaker at snow/ice regions. The decline is likely explained by the lack of SWIR

598    bands and the small thermal contrast between the clouds and the surface during the summer

27

599    nighttime in polar regions. In this case, the ML-based algorithms are not able to compensate for

600    the missing physical signatures.

601    For the daytime cloud thermodynamic phase comparison, we showed that the two RF models

602    are comparable with the MODIS MYD06 OP-Phase product, and are among the top 3 phase

603    algorithms for all surface types. The MODIS MYD06 IR-Phase, VIIRS/MODIS CLDPROP OP-

604    Phase, and CT-Phase have either relatively lower TPRs or higher FPRs over certain surface types,

605    such as shrubland, snow/ice, and barren regions. For nighttime clouds, the RF all-day model works

606    better than both CLDPROP CT-Phase and MYD06 IR-Phase for all surface types.

607    In this study, we have demonstrated the advantages of using ML-based (specifically, RF)

608    models in cloud masking and thermodynamic phase detection. In contrast with hand-tuned

609    methods, the RF models can be efficiently trained and tested for different surface types and using

610    different input variables. Meanwhile, for aerosol-free, homogeneous samples, the two RF models

611    show better and more consistent performance over different regions and surface types in

612    comparison with existing VIIRS and MODIS datasets. For more complicated scenes, RF

613    probabilities are more informative than binary mask/phase designations. However, further

614    investigation is required to understand how to use probabilities more quantitatively.

615    In the future, more spectral bands and/or spatial patterns can be used to improve pixel

616    classification skills, such as including more pixel types (e.g., dust and smoke). It is convenient to

617    apply RF models or other similar ML-based models to other instruments similar to VIIRS with the

618    help of active instruments. Most importantly, cloud mask and thermodynamic phase products from

619    well-trained RF models could be used to train other instruments in the absence of active sensors.

620    For example, the current RF model based VIIRS cloud mask/phase data could be used as reference

621    to train ML-based models for other instruments, such as MODIS, ABI/AHI, SEVIRI, and airborne

622    instruments. It remains as future work to determine how such an approach might lead to improved

623    consistency in cloud properties derived from different satellite imagers.

624        It is also important to emphasize that the model performance is highly reliant on the quality of

625    the training samples and reference labels. For example, in this study, more than 98% of the training

626    data have a VZA less than 50°, leading to more uncertain cloud phase fractions at large VZAs.

627    Using synthetic training data generated with reliable radiative transfer models could be a possible

628    way to mitigate this artifact.

642

643

644 **Reference:**

645 Ackerman, S. A., Holz, R. E., Frey, R., Eloranta, E. W., Maddux, B. C., and McGill, M., Cloud
646     detection with MODIS. Part II: Validation, *J. Atmos. Oceanic Technol.*, **25,** 1073–1086, doi:
647     10.1175/2007JTECHA1053.1, 2008.

648 Ackerman, S. A., Frey, R., Heidinger, A., Li, Y., Walther, A., Platnick, S., Meyer, K., Wind, G.,
649     Amarasinghe, N., Wang, C., Marchant, B., Holz, R. E., Dutcher, S., Hubanks, P., EOS MODIS
650     and SNPP VIIRS Cloud Properties: User guide for climate data record continuity Level-2 cloud
651     top and optical properties product (CLDPROP), version 1, 2019.

652 Baum, B. A., Menzel, W. P., Frey, R. A., Tobin, D. C., Holz, R. E., Ackerman, S. A., Heidinger,
653     A. K., and Yang, P., MODIS cloud-top property refinements for Collection 6, *J. Appl. Meteor.*
654     *Climatol.*, **51,** 1145-1163, doi: 10.1175/JAMC-D-11-0203.1, 2012.

655 Breiman, L., Random forests - random features. Technical report, University of California at
656     Berkeley, Berkeley, California, 1999.

657 Brodzik M. J., and Stewart J. S., Near-Real-Time SSM/I-SSMIS EASE-Grid Daily Global Ice
658     Concentration and Snow Extent, Version 5, doi:10.5067/3KB2JPLFPK3R, 2016.

659 Cao, C., Xiong, J., Blonski, S., Liu, Q., Uprety, S., Shao, X., Bai, Y., and Weng, F., Suomi NPP
660     VIIRS sensor data record verification, validation, and long-term performance monitoring, *J.*
661     *Geophys. Res. Atmos.*, **118,** 11,664-11,678, doi:10.1002/2013JD020418, 2013.

662 Cho, H., Nasiri, S. L., and Yang, P., Application of CALIOP Measurements to the Evaluation of
663     Cloud Phase Derived from MODIS Infrared Channels, *J. Appl. Meteor. Climatol.*, **48,** 2169-
664     2180, doi:10.1175/2009JAMC2238.1, 2009.

665 Dietterich, T. G., Ensemble methods in machine learning. International Workshop on Multiple
666     Classifier Systems, MCS 2000, Lecture Notes in Computer Science, **vol. 1857**, Springer,
667     Berlin, Heidelberg, 2000.

668 Freund, Y., An Adaptive Version of the Boost by Majority Algorithm, in Machine Learning, **43,**
669     293-318, 2001.

670 Frey, R. A., Ackerman, S. A., Liu, Y., Strabala, K. I., Zhang, H., Key, J. R., and Wang, X.: Cloud
671     detection with MODIS. Part I: Improvements in the MODIS cloud mask for Collection 5, *J.*
672     *Atmos. Oceanic Technol.,* **25,** 1057–1072, doi:10.1175/2008JTECHA1052.1, 2008.

673 Friedman, J. H., Greedy function approximation: a gradient boosting machine, *Ann. Stat.*, **29,**
674     1189–1232, 2001.

675 Gelaro, R., et al., The Modern-Era Retrospective Analysis for Research and Applications, Version
676     2 (MERRA-2), *J. Climate*, **30,** 5419–5454, doi:10.1175/JCLI-D-16-0758.1, 2017.

677 Hall, D. K., and Riggs, G. A., MODIS/Aqua Snow Cover Daily L3 Global 500m SIN Grid, Version
678     6. Boulder, Colorado USA. NASA National Snow and Ice Data Center Distributed Active
679     Archive Center, doi:10.5067/MODIS/MYD10A1.006, 2016.

680 Haynes, J. M., Noh, Y. J., Miller, S. D., Heidinger, A., and Forsythe, J. M., Cloud geometric
681     thickness and improved cloud boundary detection with GEOS ABI, 15[th] Annual Symposium

on New Generation Operational Environment Satellite Systems, Phoenix, AZ, 6 - 10 January, 2019.

Heidinger, A. K., Evan, A. T., Foster, M. J., and Walther, A., A naive bayesian cloud-detection scheme derived from CALIPSO and applied within PATMOS-x, *J. Appl. Meteor. Climatol.*, **51**, 1129–1144, doi:10.1175/JAMC-D-11-02.1, 2012.

Ho, T. K, The random subspace method for constructing decision forests, *IEEE Trans. Pattern Anal. Mach. Intell*. **20,** 832–844, 1998.

Holz, R. E., Ackerman, S. A., Nagle, F. W., Frey, R., Dutcher, S., Kuehn, R. E., Vaughan, M. A., and Baum, B., Global Moderate Resolution Imaging Spectroradiometer (MODIS) cloud detection and height evaluation using CALIOP, *J. Geophys. Res.*, **113,** D00A19, doi:10.1029/2008JD009837, 2008.

Hu, X. F., Belle, J. H., Meng, X., Wildani, A., Waller, L. A., Strickland, M. J., and Liu, Y., Estimating PM2.5 concentrations in the conterminous United States using the random forest approach, *Environmental Science & Technology,* **51,** 6936–6944, doi:10.1021/acs.est.7b01210, 2017.

Ji, C. and Ma, S., Combinations of weak classifiers, *IEEE Transactions on Neural Networks*, **8**, 32–42, 1997.

Joachims, T., Text categorization with support vector machines: Learning with many relevant features. In Proceedings of the 10th European Conference on Machine Learning, 137–142, Springer-Verlag, 1998.

Justice C. O., Vermote, E., Privette J., and Sei, A., The Evolution of U.S. Moderate Resolution Optical Land Remote Sensing from AVHRR to VIIRS. Land Remote Sensing and Global Environmental Change, B. Ramachandran, C. Justice, and M. Abrams, Eds., Remote Sensing and Digital Image Processing, **11,** Springer, New York, NY., 781-806, 2011.

Kox, S., Bugliaro, L., and Ostler, A.: Retrieval of cirrus cloud optical thickness and top altitude from geostationary remote sensing, *Atmos. Meas. Tech.*, **7,** 3233–3246, doi:10.5194/amt-7-3233-2014, 2014.

Latinne, P., Debeir, O., Decaestecker, C., Limiting the number of trees in random forests, in Multiple Classifier Systems, Manchester, U.K. IEEE, **2013**, 178-187, 2001.

Lee, T. E., Miller, S. D., Turk, F. J., Schueler, C., Julian, R., Deyo, S., Dills, P., and Wang, S., The NPOESS VIIRS Day/Night Visible Sensor, *Bull. Amer. Meteor. Soc.*, **87,** 191–200, https://doi.org/10.1175/BAMS-87-2-191, 2006.

Levy, R. C., Mattoo, S., Munchak, L. A., Remer, L. A., Sayer, A. M., Patadia, F., and Hsu, N. C., The Collection 6 MODIS aerosol products over land and ocean, *Atmos. Meas. Tech.*, **6,** 2989–3034, doi:10.5194/amt-6-2989-2013, 2013.

Liu, Y., Ackerman, S. A., Maddux, B. C., Key, J. R., and Frey, R. A., Errors in cloud detection over the Arctic using a satellite imager and implications for observing feedback mechanisms, *J. Climate*, **23,** 1894–1907, doi:10.1175/2009JCLI3386.1, 2010.

Maddux, B. C., Ackerman, S. A., and Platnick, S., Viewing geometry dependencies in MODIS cloud products, *J. Atmos. Oceanic Technol.*, **27,** 1519–1528, doi:10.1175/2010JTECHA1432.1, 2010.

Martins, J. V., Tanré, D., Remer, L., Kaufman, Y., Mattoo, S., and Levy, R., MODIS cloud screening for remote sensing of aerosols over oceans using spatial variability, *Geophys. Res. Lett.*, **29**, doi:10.1029/2001GL013252, 2002.

Marchant, B., Platnick, S., Meyer, K. G., Arnold, G. T., and Riedi, J., MODIS Collection 6 shortwave-derived cloud phase classification algorithm and comparisons with CALIOP, *Atmos. Meas. Tech.*, **9,** 1587–1599, doi:10.5194/amt-9-1587-2016, 2016.

Marchant, B., Platnick, S., Meyer, K., and Wind, G.: Evaluation of the Aqua MODIS Collection 6.1 multilayer cloud detection algorithm through comparisons with CloudSat CPR and CALIPSO CALIOP products, *Atmos. Meas. Tech. Discuss.*, doi:10.5194/amt-2019-448, in review, 2020.

McGill, M. J., Yorks, J. E., Scott, V. S., Kupchock, A. W., and Selmer, P. A., The Cloud-Aerosol Transport System (CATS): A technology demonstration on the *International Space Station, Proc. SPIE* **9612**, Lidar Remote Sensing for Environmental Monitoring XV, 96120A, doi:10.1117/12.2190841, 2015.

Meyer, K. G., Platnick, S., Arnold, G. T., Holz, R. E., Veglio, P., Yorks, J. E., and Wang, C., Cirrus cloud optical and microphysical property retrievals from eMAS during SEAC4RS using bi-spectral reflectance measurements within the 1.88 µm water vapor absorption band, *Atmospheric Measurement Techniques*, **9** (4), 1743-1753, doi:10.5194/amt-9-1743-2016, 2016.

Noel, V., Chepfer, H., Chiriaco, M., and Yorks, J.: The diurnal cycle of cloud profiles over land and ocean between 51° S and 51° N, seen by the CATS spaceborne lidar from the International Space Station, *Atmos. Chem. Phys.*, **18,** 9457–9473, doi:10.5194/acp-18-9457-2018, 2018.

Oshiro T. M., Perez P. S., Baranauskas J. A., How many trees in a random forest, in Machine Learning and Data Mining in Pattern Recognition. MLDM 2012. Lecture Notes in Computer Science, **7376,** Springer, Berlin, Heidelberg, 2012.

Pavolonis, M. J., Heidinger, A. K., and Uttal, T., Daytime global cloud typing from AVHRR and VIIRS: Algorithm description, validation, and comparisons, *J. Appl. Meteor.*, **44,** 804–826, 2005.

Pedregosa, F. et al., Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* **12,** 2825–2830, 2011.

Platnick, S., Meyer, K. G., King, M. D., Wind, G., Amarasinghe, N., Marchant, B., Arnold, G. T., Zhang, Z., Hubanks, P. A., Holz, R. E., Yang, P., Ridgway, W. L., Riedi, J.: The MODIS cloud optical and microphysical products: Collection 6 updates and examples from Terra and Aqua, *IEEE Transactions on Geoscience and Remote Sensing*, **55,** 502-525, doi: 10.1109/TGRS.2016.2610522, 2017.

Remer, L. A., Kaufman, Y. J., Tanré, D., Mattoo, S., Chu, D. A., Martins, J. V., Li, R., Ichoku, C., Levy, R. C., Kleidman, R. G., Eck, T. F., Vermote, E., and Holben, B. N., The MODIS aerosol algorithm, products, and validation, *J. Atmos. Sci.*, **62**, 947-973, doi:10.1175/JAS3385.1, 2005.

Sassen, K., and Cho, B. S., Subvisual-thin cirrus lidar dataset for satellite verification and climatological research, *American Meteorological Society*, **31,** 1275–1285. http://doi.org/10.1175/1520-0450(1992)031<1275:STCLDF>2.0.CO;2, 1992.

Sayer, A. M., Munchak, L. A., Hsu, N. C., Levy, R. C., Bettenhausen, C., and Jeong, M.-J., MODIS Collection 6 aerosol products: Comparison between Aqua's e-Deep Blue, Dark Target, and "merged" data sets, and usage recommendations, *J. Geophys. Res. Atmos.*, **119,** 13,965-13,989, doi:10.1002/2014JD022453, 2014.

Sayer, A. M., Hsu, N. C., Lee, J., Bettenhausen, C., Kim, W. V., and Smirnov, A., Satellite Ocean Aerosol Retrieval (SOAR) algorithm extension to S-NPP VIIRS as part of the "Deep Blue" aerosol project, *J. Geophys. Res. Atmos.*, **123,** doi:10.1002/2017JD027412, 2017.

Scornet, E., Tuning parameters in random forests. ESAIM: Procs, 60: 144–162, 2018.

Seemann, S. W., Borbas, E. E., Knuteson, R. O., Stephenson, G. R., and Huang, H., Development of a global infrared land surface emissivity database for application to clear sky sounding retrievals from multispectral satellite radiance measurements, *J. Appl. Meteor. Climatol.*, **47,** 108–123, 2008.

Stephens, G. L., et al., The CloudSat mission and the A-Train: A new dimension of space-based observations of clouds and precipitation, *Bull. Amer. Meteorol. Soc.*, **83,** 1771-1790, doi:10.1175/BAMS-83-12-1771, 2002.

Strandgren, J., Bugliaro, L., Sehnke, F., and Schröder, L.: Cirrus cloud retrieval with MSG/SEVIRI using artificial neural networks, *Atmos. Meas. Tech.*, **10,** 3547–3573, doi:10.5194/amt-10-3547-2017, 2017.

Stubenrauch, C. J., Rossow, W. B., Kinne, S., Ackerman, S., Cesana, G., Chepfer, H., Di Girolamo, L., Getzewich, B., Guignard, A., Heidinger, A., Maddux, B. C., Menzel, W. P., Minnis, P., Pearl, C., Platnick, S., Poulsen, C., Riedi, J., Sun-Mack, S., Walther, A., Winker, D., Zeng, S., and Zhao, G., Assessment of Global Cloud Datasets from Satellites: Project and Database Initiated by the GEWEX Radiation Panel, *Bull. Amer. Meteor. Soc.*, **94,** 1031–1049, doi:10.1175/BAMS-D-12-00117.1, 2013.

Sulla-Menashe, D., and Friedl, M. A., User Guide to Collection 6 MODIS Land Cover (MCD12Q1 and MCD12C1) Product; USGS: Reston, VA, USA, 2018.

Tanelli, S., Durden, S. L., Im, E., Pak, K., Reinke, D., Partain, P., Haynes, J., and Marchand, R., CloudSat's cloud profiling radar after two years in orbit: Performance, calibration, and processing, *IEEE Trans. Geosci. Remote Sens.*, **46,** 3560–3573, doi:10.1109/TGRS.2008.2002030, 2008.

Thampi, B. V., Wong, T., Lukashin, C., and Loeb, N. G., Determination of CERES TOA fluxes using machine learning algorithms. Part I: Classification and retrieval of CERES cloudy and clear scenes, *J. Atmos. Oceanic Technol.*, **34,** 2329–2345, doi:10.1175/JTECH-D-16-0183.1, 2017.

Tumer, K., and Ghosh, J., Error correlation and error reduction in ensemble classifiers, *Connection Science,* **8**, 385-403, doi:10.1080/095400996116839, 1996.

Wan, Z., Zhang, Y., Zhang, Q., and Li, Z.-L., Quality assessment and validation of the MODIS global land surface temperature, *Int. J. Remote Sens.*, **25**, 261–274, doi:10.1080/0143116031000116417, 2004.

803 Wang, C., Yang, P., Dessler, A., Baum, B. A., and Hu, Y., Estimation of the cirrus cloud scattering
804     phase function from satellite observations, *Journal of Quantitative Spectroscopy and Radiative*
805     *Transfer*, **138**, 36-49 doi:10.1016/j.jqsrt.2014.02.001, 2014.

806 Wang, C., Platnick, S., Zhang, Z., Meyer, K., and Yang, P., Retrieval of ice cloud properties using
807     an optimal estimation algorithm and MODIS infrared observations: 1. Forward model, error
808     analysis, and information content, *J. Geophys. Res. Atmos.*, **121,** 5809-5826
809     doi:10.1002/2015jd024526, 2016a.

810 Wang, C., Platnick, S., Zhang, Z., Meyer, K., Wind, G., and Yang, P., Retrieval of ice cloud
811     properties using an optimal estimation algorithm and MODIS infrared observations: 2.
812     Retrieval evaluation, *J. Geophys. Res. Atmos.*, **121**, doi:10.1002/2015jd024528, 2016b.

813 Wang, C., Platnick, S., Fauchez, T., Meyer, K., Zhang, Z., Iwabuchi, H., and Kahn, B. H., An
814     assessment of the impacts of cloud vertical heterogeneity on global ice cloud data records from
815     passive satellite retrievals, *Journal of Geophysical Research: Atmospheres*, **124,** 1578-1595.
816     doi:10.1029/2018JD029681, 2019.

817 Winker, D. M., Tackett, J. L., Getzewich, B. J., Liu, Z., Vaughan, M. A., and Rogers, R. R., The
818     global 3-D distribution of tropospheric aerosols as characterized by CALIOP, *Atmos. Chem.*
819     *Phys.*, **13,** 3345-3361, doi:10.5194/acp-13-3345-2013, 2013.

820 Wolters, E. L., Roebeling, R. A., and Feijt, A. J., Evaluation of cloud-phase retrieval methods for
821     SEVIRI on Meteosat-8 using ground-based lidar and cloud radar data, *J. Appl. Meteor.*
822     *Climatol.*, **47,** 1723–1738, doi:10.1175/2007JAMC1591.1, 2008.

823 Wu, Y., de Graaf, M., and Menenti, M., Improved MODIS Dark Target aerosol optical depth
824     algorithm over land: angular effect correction, *Atmos. Meas. Tech.*, **9,** 5575-5589,
825     doi:10.5194/amt-9-5575-2016, 2016.

826 Yuan, T., Wang, C., Song, H., Platnick, S., Meyer, K., and Oreopoulos, L., Automatically finding
827     ship tracks to enable large-scale analysis of aerosol-cloud interactions, *Geophysical Research*
828     *Letters*, **46,** 7726– 7733, doi: 10.1029/2019GL083441, 2019.

829

830

831

832

833

834

835

836

837

838

839

840

841    Table 1. Existing VIIRS and MODIS cloud mask and phase products used for comparison. Note
842    that MYD35 and MYD06 are the standard MODIS Aqua products, and CLDMSK and CLDPROP
843    are the MODIS Aqua and VIIRS common algorithm continuity products.

844

| Instrument | Cloud Mask | Cloud Phase |
|---|---|---|
| MODIS | MYD35 V6.1 | MYD06 IR-Phase V6.1 |
| | | MYD06 OP-Phase V6.1 |
| | CLDMSK V1.0 | CLDPROP CT-Phase V1.0 |
| | | CLDPROP OP-Phase V1.1 |
| VIIRS | CLDMSK V1.0 | CLDPROP CT-Phase V1.0 |
| | | CLDPROP OP-Phase V1.1 |

845

846

847     Table 2: Data collection strategies and the number of pixels for all surface types.
848

| # of VIIRS 750m pixels (million) | Condition | Ocean | Forest | Cropland | Grass | Barren | Shrub | Snow/Ice | Total |
|---|---|---|---|---|---|---|---|---|---|
| All collocation | None | 219.7 | 18.7 | 8.7 | 17.5 | 17.1 | 13.6 | 37.4 | 332.7 |
| Aerosol Free | CALIOP Aerosol 5km column AOD < 0.05 | 142.6 | 13.0 | 3.7 | 10.0 | 10.5 | 9.3 | 34.3 | 223.2 |
| Clear | Aerosol Free, Cloud 1km Layer = 0 | 17.7 | 2.5 | 1.5 | 1.8 | 2.9 | 3.1 | 13.1 | 42.5 |
| **Clear (homogeneous)** | **Aerosol Free, Cloud 1km/5km Layer = 0** | **15.2** | **2.3** | **1.5** | **1.7** | **2.7** | **3.0** | **12.7** | **39.1** |
| Cloudy | Aerosol Free, Cloud 1km Layer > 0 | 124.9 | 10.5 | 2.1 | 8.1 | 7.7 | 6.2 | 21.2 | 180.7 |
| Cloudy (homogeneous) | Aerosol Free, Cloud 1km/5km Layer > 0 | 115.5 | 9.5 | 1.8 | 7.4 | 6.6 | 5.3 | 15.8 | 162.0 |
| Single Phase Cloud | Aerosol Free, Cloud 1km Liquid or Ice Phase | 65.1 | 4.4 | 1.0 | 4.0 | 3.4 | 2.4 | 13.5 | 93.7 |
| **Single Phase Cloud (homogeneous)** | **Aerosol Free, Cloud 1km/5km Liquid or Ice Phase** | **64.2** | **4.3** | **0.9** | **3.9** | **3.3** | **2.3** | **12.7** | **91.5** |
| **Liquid Phase Cloud (homogeneous)** | **Aerosol Free, Cloud 1km/5km Liquid Phase** | **40.5** | **1.8** | **0.3** | **1.7** | **1.3** | **1.0** | **3.2** | **49.7** |
| **Ice Phase Cloud (homogeneous)** | **Aerosol Free, Cloud 1km/5km Ice Phase** | **23.7** | **2.5** | **0.6** | **2.2** | **2.0** | **1.3** | **9.5** | **41.8** |

849
850

851 Table 3: Accuracy scores of RF all-day models based on testing pixels with different inputs and a
852 fixed model configuration (N_Trees = 150 and Max_TreeDepths = 15).

| #<br>Input | Model Input | Ocean | Forest | Shrubland | Crop | Grassland | Barren | Snow/Ice | All<br>Surface* |
|---|---|---|---|---|---|---|---|---|---|
| 1 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, and VZA | 90.3 | 89.9 | 88.7 | 88.4 | 88.2 | 88.0 | 87.4 | 89.4 |
| 2 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, VZA, and Lat/Lon | 92.1 | 90.1 | 89.8 | 90.7 | 89.5 | 90.1 | 88.0 | 90.9 |
| 3 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, VZA, and $T_S$ | 93.1 | 90.9 | 89.9 | 91.4 | 90.2 | 90.3 | 88.5 | 91.7 |
| 4 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, VZA, Lat/Lon, and $T_S$ | 93.2 | 91.7 | 90.0 | 91.8 | 91.2 | 90.8 | 88.9 | 92.0 |
| 5 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, VZA, $T_S$, and $\varepsilon_S$ | 93.2 | 91.4 | 89.8 | 91.4 | 90.4 | 90.4 | 88.8 | 91.9 |
| 6 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, VZA, Lat/Lon, $T_S$, and $\varepsilon_S$ | 93.2 | 91.8 | 90.1 | 91.8 | 91.3 | 90.6 | 88.9 | 92.0 |

853 *The all-surface accuracy scores are weighted by pixel numbers of individual surface types.

854  Table 4: Accuracy scores of RF daytime models based on testing pixels with different inputs and
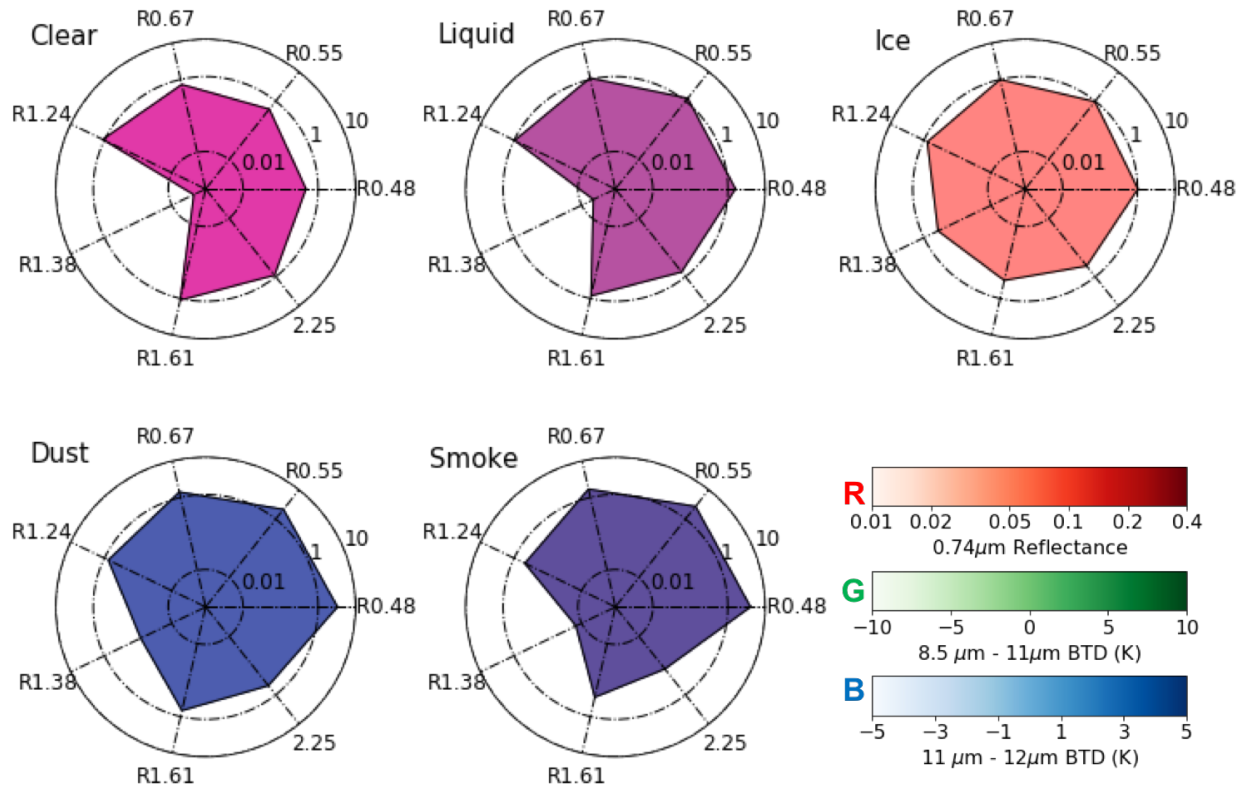855  a fixed model configuration (N_Trees = 150 and Max_TreeDepths = 15).

| # Input | Model Input | Ocean | Forest | Shrubland | Crop | Grassland | Barren | Snow/Ice | All Surface* |
|---|---|---|---|---|---|---|---|---|---|
| 1 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, $R_{0.86}$, $R_{1.38}$, $R_{1.61}$, $R_{2.25}$, VZA, and SZA | 95.47 | 93.71 | 93.25 | 93.86 | 92.82 | 94.04 | 94.94 | 94.97 |
| 2 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, $R_{0.86}$, $R_{1.38}$, $R_{1.61}$, $R_{2.25}$, VZA, SZA, and RAA | 95.47 | 93.72 | 93.22 | 93.84 | 92.81 | 94.02 | 94.94 | 94.97 |
| 3 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, $R_{0.86}$, $R_{1.38}$, $R_{1.61}$, $R_{2.25}$, Lat/Lon, VZA, and SZA | 95.47 | 93.74 | 93.36 | 93.95 | 92.95 | 94.16 | 94.95 | 94.99 |
| 4 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, $R_{0.86}$, $R_{1.38}$, $R_{1.61}$, $R_{2.25}$, $R_{1.24}$, Lat/Lon, VZA and SZA | 95.51 | 93.73 | 93.47 | 93.93 | 92.98 | 94.21 | 95.05 | 95.04 |
| 5 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, $R_{0.86}$, $R_{1.38}$, $R_{1.61}$, $R_{2.25}$, Ts, Lat/Lon, VZA, SZA, and RAA | 95.45 | 93.77 | 93.36 | 93.93 | 92.92 | 94.21 | 94.95 | 94.98 |
| 6 | $BT_{8.6}$, $BT_{11}$, $BT_{12}$, $R_{0.86}$, $R_{1.38}$, $R_{1.61}$, $R_{2.25}$, $R_{0.48}$, $R_{0.67}$, $R_{1.24}$, VZA, and SZA | 95.51 | 93.90 | 93.54 | 94.11 | 93.07 | 94.38 | 95.17 | 95.09 |

856  *The all-surface accuracy scores are weighted by pixel numbers of individual surface types.

857 Table 5: Fractions of the 2017 validation samples that have determined phases (i.e., liquid water
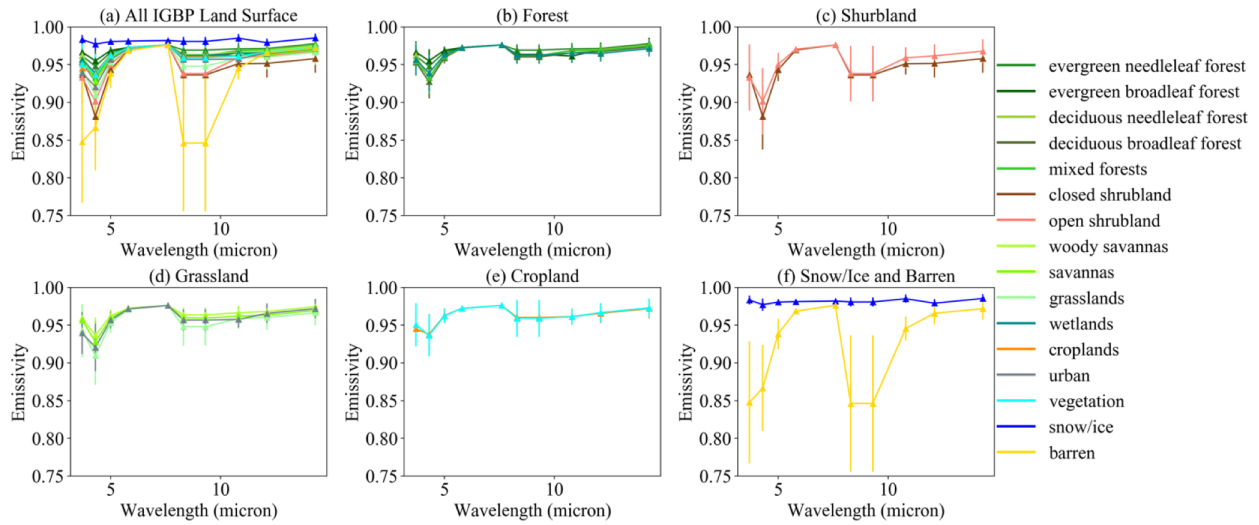858 or ice) in different surface types.
859

| Determined Phase (%) | Ocean | Forest | Shrubland | Crop | Grassland | Barren | Snow/Ice | All |
|---|---|---|---|---|---|---|---|---|
| MODIS MYD06 IR-Phase | 89 | **75** | **74** | 80 | **79** | **75** | **66** | 85 |
| MODIS MYD06 OP-Phase | 97 | 99 | 97 | 98 | 99 | 95 | 92 | 97 |
| MODIS CLDPROP OP-Phase | 98 | 99 | 98 | 99 | 99 | 97 | 99 | 98 |
| VIIRS CLDPROP OP-Phase | 98 | 99 | 97 | 99 | 98 | 96 | 99 | 98 |

860

Figure 1. Spectral patterns of the five different pixel types (averaged over 1,000 pixels for each type). For each plot, an apex indicates reflectance ratio between a given VNIR/SWIR band and the 0.86-$\mu$m band, and the spread is filled by false RGB composite (Red: 0.74-$\mu$m reflectance; Green: 8.5-11$\mu$m brightness temperature difference (BTD); Blue: 11-12$\mu$m BTD). The spectral patterns are used in the machine learning algorithms.
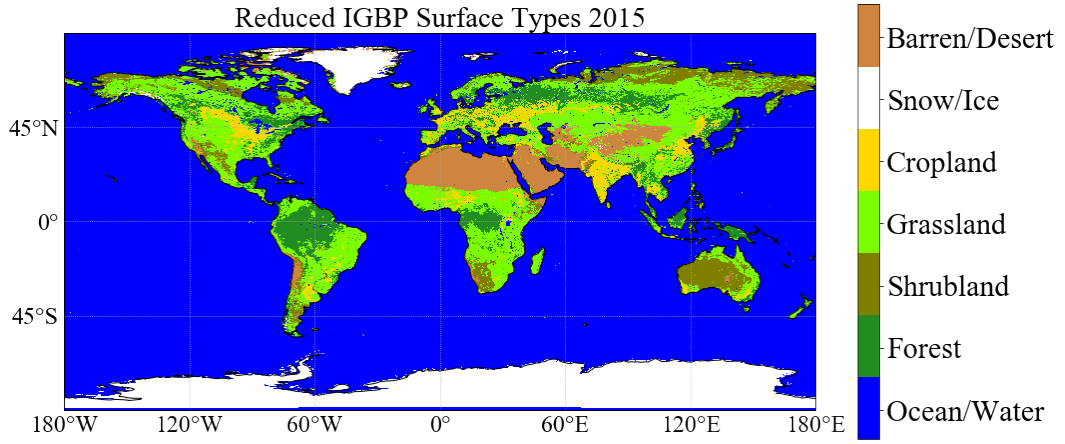
868

Figure 2. Climatology of the spectral surface emissivity data from the UW-Madison baseline fit land surface emissivity database [*Seemann et al.*, 2008] for different IGBP surface types. Error bars indicate the emissivity standard deviations at given wavelengths.

872

873
874 Figure 3. Climatology of the spectral surface white sky surface albedo data from MCD12C1 [*Sulla-*
875 *Menashe and Friedl* 2018] for different IGBP surface types. Error bars indicate the albedo standard
876 deviations at given wavelengths.
877

Reduced IGBP Surface Types 2015

Barren/Desert
Snow/Ice
Cropland
Grassland
Shrubland
Forest
Ocean/Water

878

879    Figure 4. A global map of the seven reduced surface types chosen for the RF model training.

880

Figure 5. Global distributions of the of clear and cloudy pixels from collocated VIIRS and CALIOP data from 2013 to 2017. Panels a) and d) show the total clear and cloudy pixel counts, respectively. Panels b) and d) show the pixel counts after applying the quality control. The corresponding selection ratios are shown in panels c) and f).
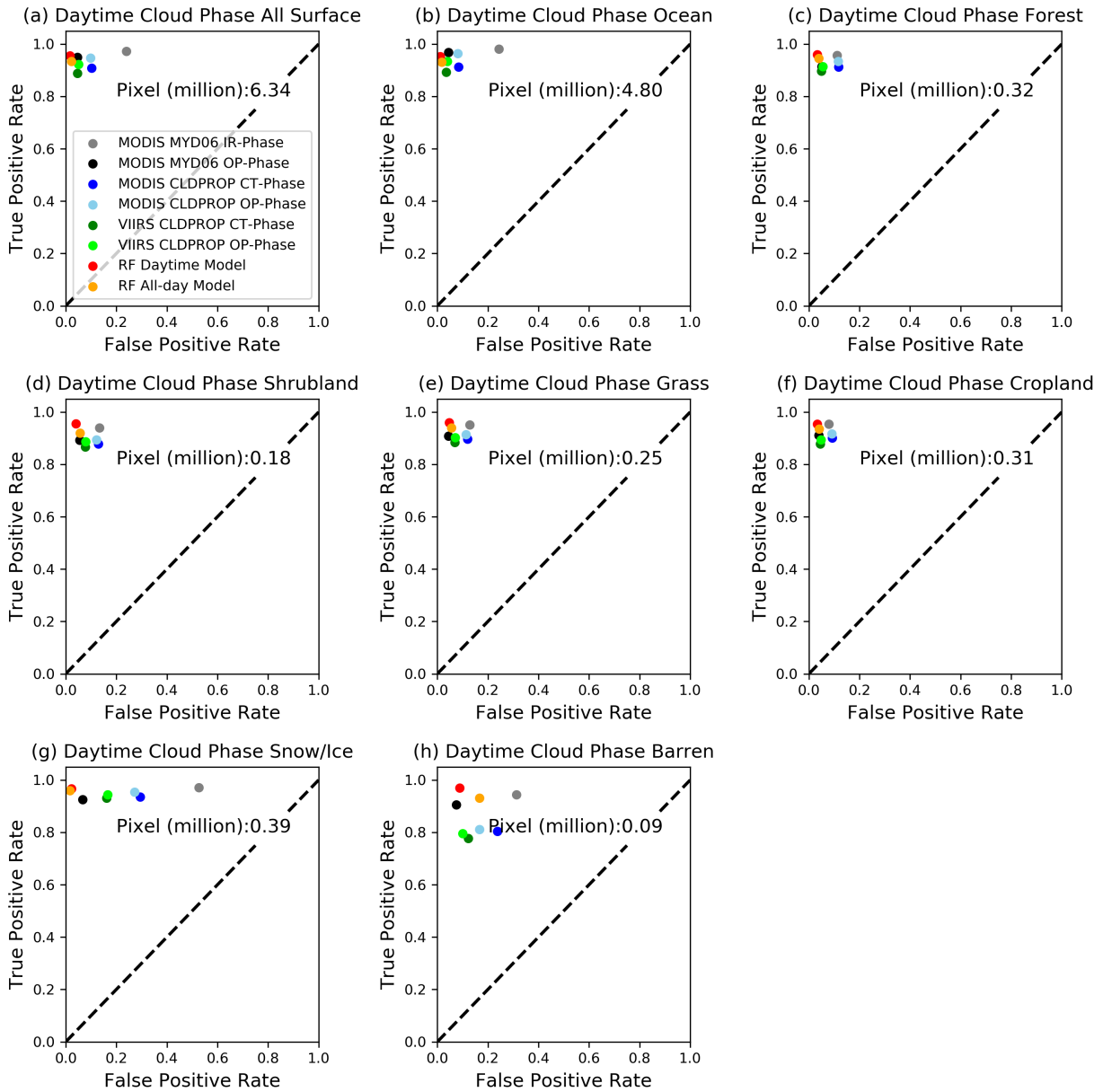
Figure 6. False Positive Rate (FPR) versus True Positive Rate (TPR) plots of daytime cloud mask from the two RF models and operational algorithms. Collocated CALIOP Level 2 products in 2017 are used as reference. Global comparisons are shown in panel (a), while panels (b) through (h) show comparisons for difference surface types. The total pixel number is shown in each panel.

(a) Nighttime Cloud Mask All Surface

(b) Nighttime Cloud Mask Ocean

(c) Nighttime Cloud Mask Forest

(d) Nighttime Cloud Mask Shrubland

(e) Nighttime Cloud Mask Grass

(f) Nighttime Cloud Mask Cropland

(g) Nighttime Cloud Mask Snow/Ice
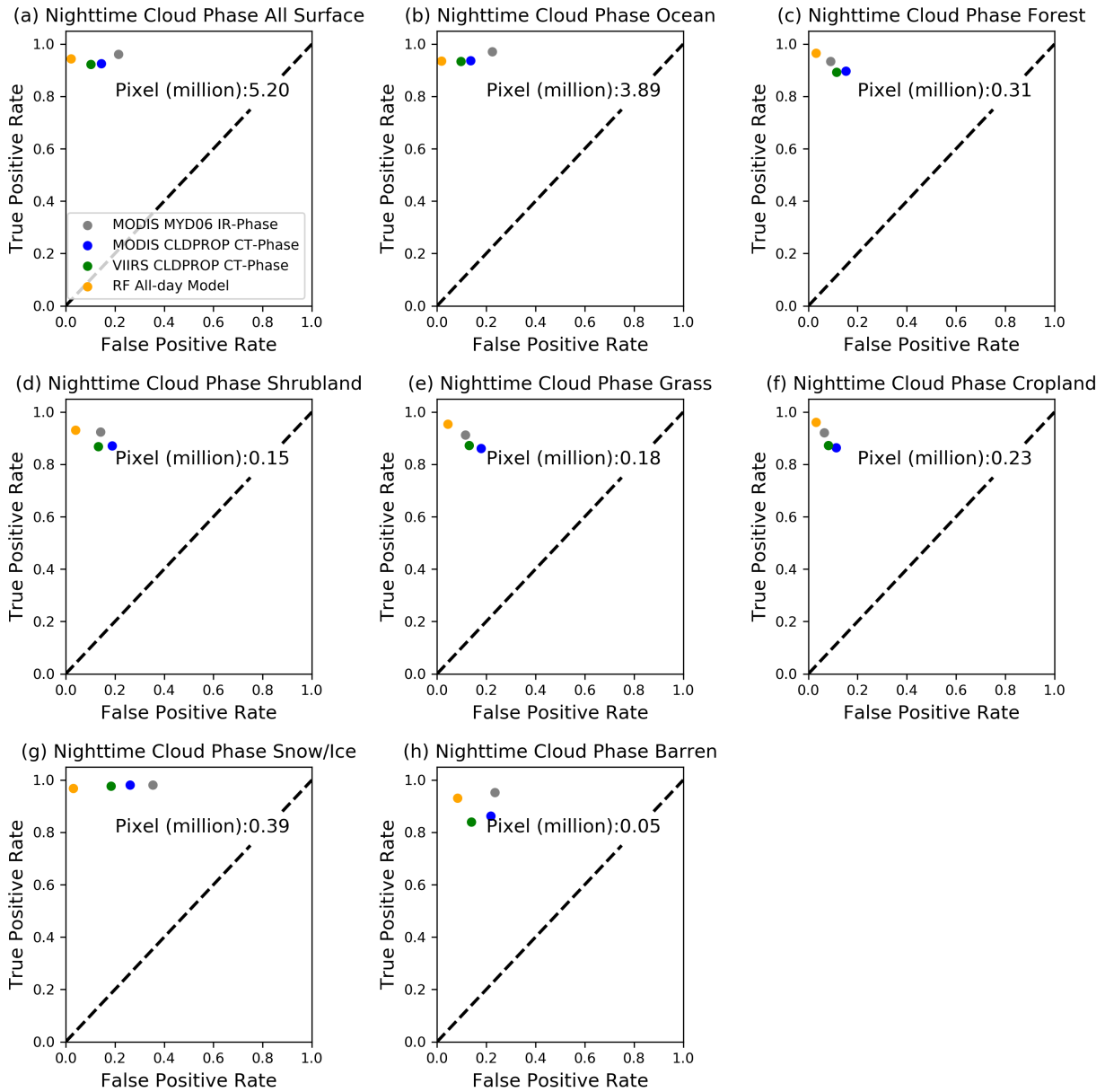
(h) Nighttime Cloud Mask Barren

893

894    Figure 7. Similar to Figure 6, but for nighttime cloud mask comparisons. The total pixel number
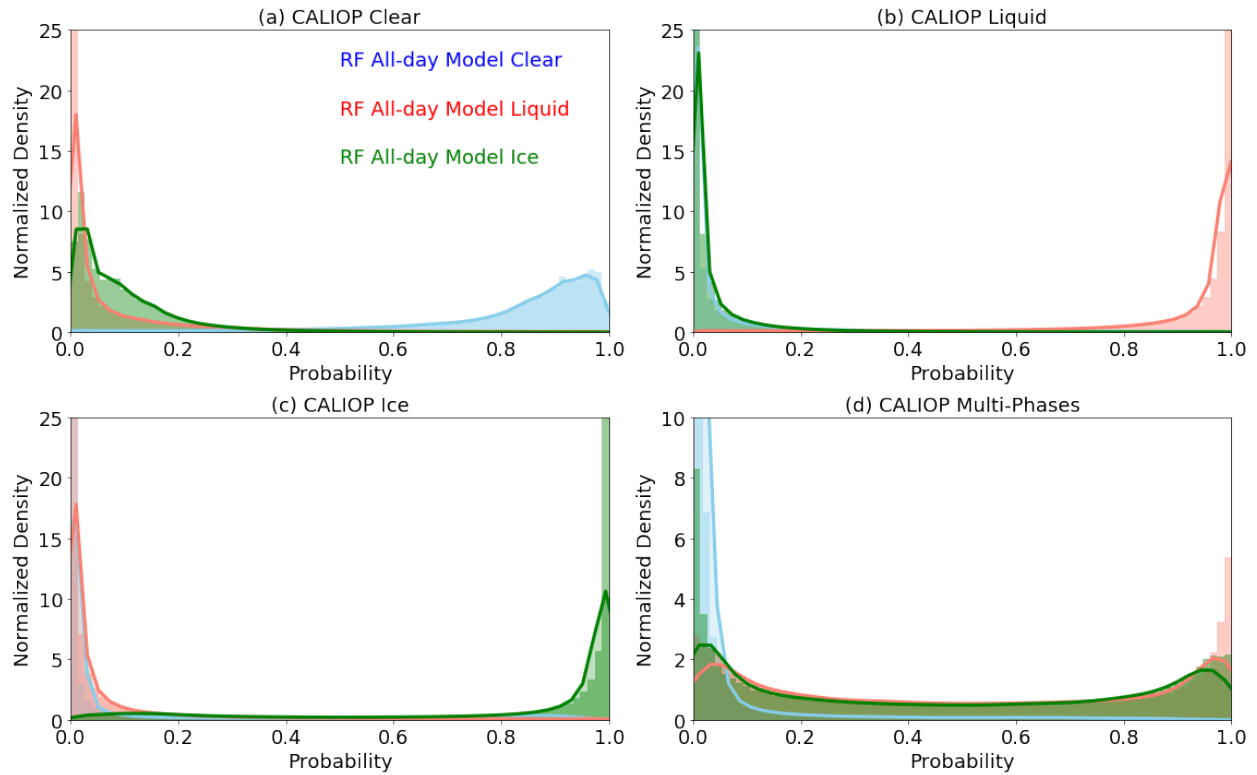895    is shown in each panel.

896

46

Figure 8. Similar to Figure 6, but for daytime cloud thermodynamic phase comparisons. The total pixel number is shown in each panel. Note that for specific products, the total pixel numbers are less because of the exclusion of "unknown phase" category (see text for more details).

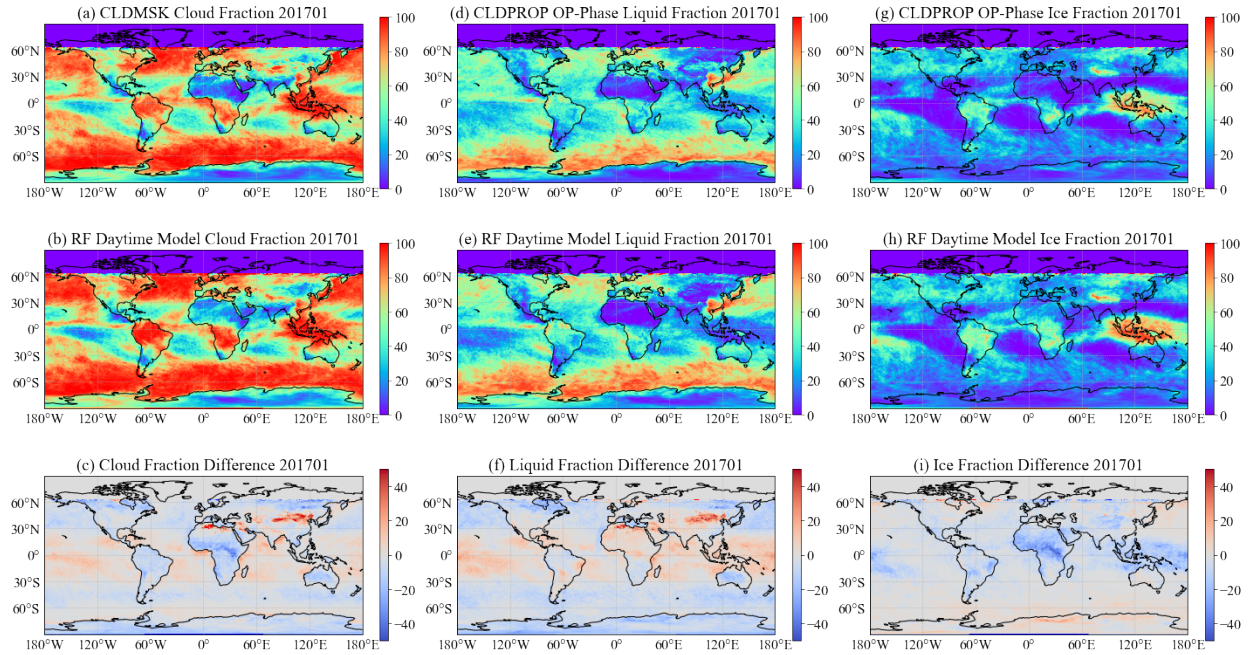(a) Nighttime Cloud Phase All Surface
(b) Nighttime Cloud Phase Ocean
(c) Nighttime Cloud Phase Forest
(d) Nighttime Cloud Phase Shrubland
(e) Nighttime Cloud Phase Grass
(f) Nighttime Cloud Phase Cropland
(g) Nighttime Cloud Phase Snow/Ice
(h) Nighttime Cloud Phase Barren

Pixel (million):5.20
Pixel (million):3.89
Pixel (million):0.31
Pixel (million):0.15
Pixel (million):0.18
Pixel (million):0.23
Pixel (million):0.39
Pixel (million):0.05

MODIS MYD06 IR-Phase
MODIS CLDPROP CT-Phase
VIIRS CLDPROP CT-Phase
RF All-day Model

902

Figure 9. Similar to Figure 6, but for nighttime cloud thermodynamic phase comparisons. The total
pixel number is shown in each panel. Note that for specific products, the total pixel numbers are
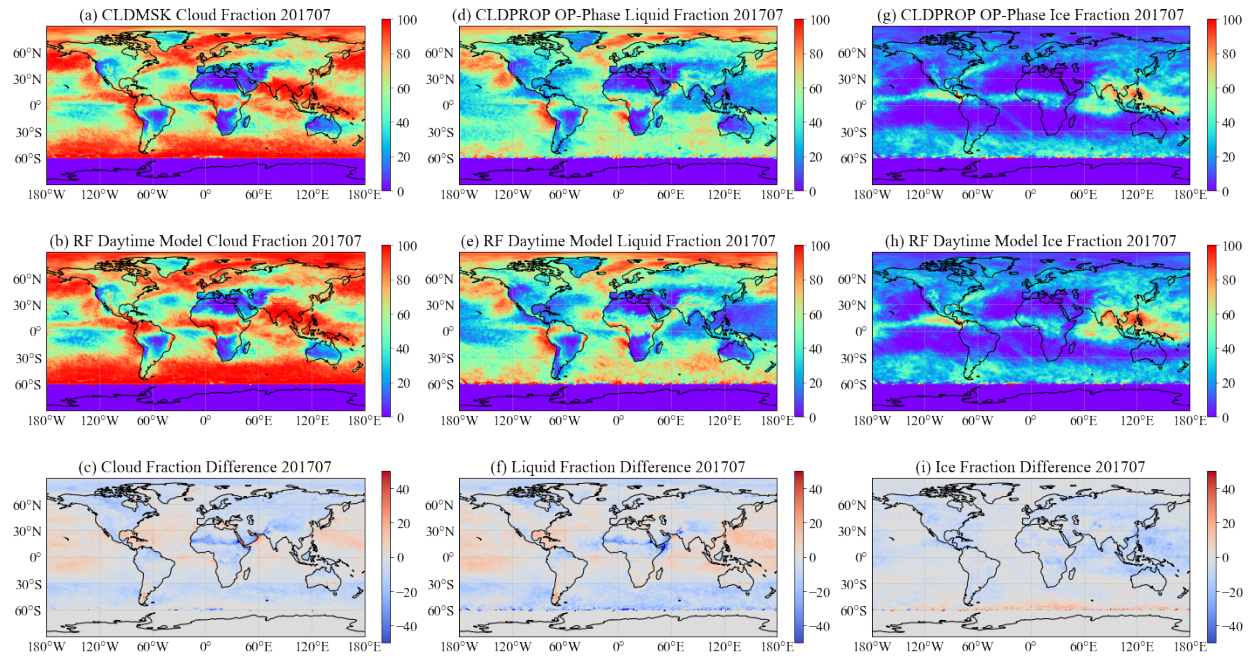less because of the exclusion of "unknown phase" category (see text for more details).

906

48

Figure 10. Normalized density functions of the clear (blue), liquid water cloud (red), and ice cloud (green) probabilities from the RF all-day model in four CALIOP detected aerosol-free scenes: (a) clear, (b) homogenous liquid, (c) homogenous ice, and (d) multi-layer cloud with different thermodynamic phases.
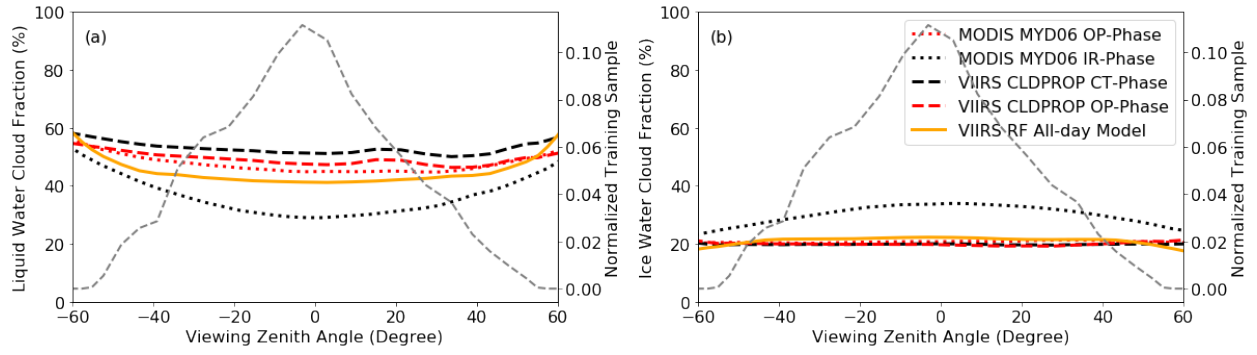
913

Figure 11. Comparisons between one-month daytime cloud mask and thermodynamic phase products from the VIIRS CLDMSK and CLDPROP OP-Phase (top row) and the RF daytime model (second row), and their differences (VIIRS – RF daytime, bottom row) in January, 2017.

914
915
916
917

918

Figure 12. Similar to Figure 11, but for comparisons in July, 2017.

920

Figure 13. Liquid water (a) and ice (b) cloud fractions as a function of viewing zenith angle from the one-month daytime cloud mask/phase products in January 2017. The gray dashed curve is the probability density function of the 4-year VIIRS/CALIOP training samples (2013-2016).