1 **Cloud-Precipitation Hybrid Regimes and their Projection onto IMERG**

2 **Precipitation Data**

3

4 Daeho Jin[a,b], Lazaros Oreopoulos[b], Dongmin Lee[c,b], Jackson Tan[a,b], Nayeong Cho[a,b]

5 [a] *Universities Space Research Association, Columbia, MD, USA*

6 [b] *NASA Goddard Space Flight Center, Greenbelt, MD, USA*

7 [c] *Morgan State University, Baltimore MD, USA*

8

9

10 *Corresponding author*: Daeho Jin, Daeho.Jin@NASA.gov

11

12                                    ABSTRACT

13      We extend and enhance the concept of the Cloud Regimes (CRs) developed from two-

14  dimensional joint histograms of cloud optical thickness and cloud top pressure from the

15  Moderate Resolution Imaging Spectroradiometer (MODIS), by adding precipitation

16  information in order to better understand cloud-precipitation relationships. Taking advantage

17  of the high-resolution Integrated Multi-satellitE Retrievals for GPM (IMERG) precipitation

18  dataset, cloud-precipitation "hybrid" regimes are derived by implementing the $k$-means

19  clustering algorithm with advanced initialization and objective measures to determine the

20  most optimal clusters. By expressing precipitation rates within 1-degree grid cell as

21  histograms and making choices on the relative weight of cloud and precipitation, we could

22  obtain several editions of hybrid cloud-precipitation regimes (CPRs), and examine their

23  characteristics.

24      In the deep tropics, when precipitation is weighted weakly, the cloud part of the hybrid

25  centroids resembles the centroid of cloud-only regimes, but still tightens the cloud-

26  precipitation relationship by decreasing the precipitation variability of each regime. As

27  precipitation weight progressively increases, the shape of the cloudy part of the hybrid

28  centroids becomes blunter, while the precipitation part of the centroids sharpens. In the case

29  where cloud and precipitation are weighted equally, the CPRs representing high clouds with

30  intermediate to heavy precipitation exhibit distinct features in the precipitation parts of the

31  centroids, which allows us to project them onto the 30-minly IMERG domain. Such a

32  projection can be used to overcome the temporal sparseness of MODIS cloud observations,

33  which leads to great application potential for various convection-focused studies, including

34  diurnal cycle analysis.

35                            SIGNIFICANCE STATEMENT

36    Clouds and precipitation are related closely, but in complex ways. In this work we

37    attempt to provide a classification of daytime cloud-precipitation co-occurrence and co-

38    variability, with emphasis in tropical regions. We achieve such a classification using $k$-means

39    clustering algorithm applied on cloud property and precipitation intensity histograms which

40    yields "hybrid" clusters. These hybrid clusters reveal more detailed features of coincident

41    daytime cloud and precipitation systems than clusters where clouds and precipitation are

42    treated separately. Moreover, the realization that precipitation features associated with high

43    and thick clouds have very distinct patterns enables hybrid cluster prediction based solely on

44    precipitation information, which has the important implication that rarer cloud observations

45    can be extended to the more frequent (including nighttime) precipitation domain.

46

## 1. Introduction

48    In many applications, a variable or combinations of variables that co-vary need to be

49    sorted into groups whose members are considered similar. One option to accomplish the

50    grouping is clustering analysis, a discipline of unsupervised machine learning. Among

51    various algorithms that perform clustering, "$k$-means" is one of the most popular options in

52    geophysical sciences due to its simplicity and efficiency in processing large volumes of data.

53    Examples of recent studies where $k$-means clustering is used are the grouping of precipitation

54    patterns to identify the South Pacific convergence zone (SPCZ; Pike and Lintner 2020),

55    analysis of geopotential height data to identify weather patterns for subseasonal forecast

56    (Robertson et al. 2020), and finding dominant modes in sea surface temperature data to

57    identify two kinds of the North Pacific Meridional Mode (NPMM; Zhao et al. 2020), etc.

58    *k*-means clustering has also been applied in the last two decades to cloud grouping. Based

59    on the gridded Level-3 2D-joint histogram of cloud top height (CTP) and cloud optical

60    thickness (COT) retrieved from the International Satellite Cloud Climatology Project

61    (ISCCP), dominant mixtures of clouds, later called "weather states", were identified in the

62    tropical western Pacific (Jakob and Tselioudis 2003), the deep tropics from 15°S to 15°N

63    (Rossow et al. 2005), extended tropics and mid-latitudes (Oreopoulos and Rossow 2011), and

64    globally (Tselioudis et al. 2013). The same methodology was extended to similar 2D-joint

65    histogram of CTP and COT retrieved from the Moderate Resolution Imaging

66    Spectroradiometer (MODIS), to obtain cloud groups referred to as "cloud regimes (CRs)"

67    (Oreopoulos et al. 2014, 2016; Jin et al. 2020).

68    Clouds and precipitation are closely related to each other, albeit in complex ways, so the

69    effort of Luo et al. (2017) to perform joint clustering of cloud and precipitation information

70    came as a natural progression in expanding clustering applications. Using the Tropical

71    Rainfall Measuring Mission (TRMM) Ku-band Precipitation Radar and the CloudSat W-band

72    Cloud-Profiling Radar, they first built 2D joint histogram of height and radar reflectivity

73    (a.k.a. H-dBZ histogram) for rather sparse coincident observations, on which they then

74    performed *k*-means clustering analysis. They also tested another expanded version of joint

75    histograms where CALIOP lidar products were added to capture optically thinner clouds, and

76    obtained a larger number of meaningful joint cloud-precipitation groups. This pioneering

77    work opened new pathways to group microphysical properties of hydrometeors by regimes

78    with data that can also resolve vertical structures. Combined cloud-precipitation analysis, but

79    without joint clustering, have also been performed within the framework of weather states or

80    CRs. But in these studies precipitation variability was a dependent variable sorted for specific

4

81   kinds of cloud mixtures as represented by the weather states or CRs (e.g., Lee et al. 2013;

82   Rossow et al. 2013; Tan et al. 2015; Tan and Oreopoulos 2019).

83      Recently, precipitation datasets have been greatly improved in terms of quality and

84   spatiotemporal coverage due to advances in algorithms such as the Integrated Multi-satellitE

85   Retrievals for GPM (IMERG) product providing precipitation rates at 0.1° every 30 minutes.

86   The combination of the IMERG precipitation and MODIS cloud products provides an

87   unprecedented opportunity to examine cloud-precipitation joint variability not possible with

88   previous generation datasets. We thus return in this study to the joint clustering concepts of

89   Luo et al. (2017) aiming once again to identify dominant mixtures of cloud and precipitation

90   patterns. While our data, Level-3 cloud and precipitation products, do not have the capability

91   to resolve vertical variability, we can perform joint clustering with much wider coverage

92   compared to the availability of the Level-2 reflectivity and backscatter. It turns out that the

93   existence of a tight coupling between clouds and precipitation in some of our "hybrid"

94   regimes allows us to take advantage of the higher temporal resolution of IMERG to greatly

95   expand the rarer cloud information suffering the limitations of sun-synchronous satellite

96   observations. We discuss this further in section 5 of this paper.

97      The remainder of the paper provides the details of data and $k$-means clustering

98   methodology (sections 2 and 3), formally presents the cloud-precipitation hybrid regimes and

99   discusses their characteristics in section 4. Section 6 summarizes the study and discusses

100  possible applications of the new dataset.

101

## 2. Data

103  *a. MODIS cloud data*

5

File generated with AMS Word template 1.0

104    Cloud properties are retrieved from the Moderate Resolution Imaging Spectroradiometer

105    (MODIS) instrument aboard the Terra and Aqua satellites. The MODIS cloud product

106    (MOD08_D3 and MYD08_D3; King et al. 2003; Platnick et al. 2003, 2017b) provides Level-

107    3 cloud observations at daily time scales with $1°\times1°$ horizontal resolution. Among various

108    variables in Level-3 products, we specifically use the ISCCP-like 2D joint histogram of cloud

109    optical thickness (COT) and cloud top pressure (CTP). The histogram is composed of cloud

110    fraction (CF) values along 7 classes of CTP and 6 classes of COT (for a total 42 histogram

111    bins), thus providing information about pixel-level cloud variability at the $1°$ scale. Since the

112    recent major version of the MODIS atmospheric datasets, known as "Collection 6" (Platnick

113    et al. 2017a), a separate histogram for "partially cloudy" (PCL) pixels is provided, flagged as

114    such by the so-called "clear-sky restoral" algorithm (Pincus et al. 2012; Zhang and Platnick

115    2011). The 2D joint histograms used in this study include the sum of the PCL and nominal

116    joint histograms, as in Jin et al. (2018, 2020). The update from Collection 6 to Collection 6.1

117    used here is relatively minor (Platnick et al. 2018).

118    *b. IMERG precipitation data*

119    The Integrated Multi-satellitE Retrievals for GPM (IMERG) data provides seamless

120    precipitation estimates at a $0.1°$ grid every half hour by unifying observations from a network

121    of partner satellites in the Global Precipitation Measurement (GPM) constellation (Huffman

122    et al. 2019a,b; Tan et al. 2019a). The most recent major update version V06 extends spatial

123    coverage to the entire globe (except over frozen surfaces at high latitudes) and the temporal

124    period back to June 2000 (the pre-GPM era of the Tropical Rainfall Measuring Mission –

125    TRMM) onwards. The IMERG product comprises three runs (Early, Late, and Final), of

126    which we use the Final run which is of best quality. We note that for this study we limit the

6

127    data period for both cloud and precipitation from June 2014 to May 2019 in order to avoid

128    potential risk of inconsistencies between the GPM and TRMM satellites.

129    *c. Spatio-temporal matching between MODIS and IMERG data*

130         The MODIS Level-3 gridded data is provided daily for each of the Terra and Aqua

131    satellites. Observations on swath paths for a large portion of the globe take place at similar

132    local time but varying Coordinated Universal Time (UTC). In order to temporally match

133    MODIS cloud data and IMERG precipitation data which are segmented by UTC, we

134    calculate the UTC of each MODIS grid cell using the assigned mean solar zenith angle in the

135    Level-3 product, and then select the temporally closest IMERG data point. The details of this

136    temporal matching method are described in Jin et al. (2018), and although in that paper the

137    precipitation data was the TRMM Multi-satellite Precipitation Analysis (TMPA), the

138    principle of the method is the same. Spatial matching is much easier: for each $1°\times1°$ grid cell

139    of MODIS clouds, the one hundred enclosed precipitation rates of $0.1°\times0.1°$ resolution are

140    assigned. Hence, we ultimately obtain 42 values of binned cloud fraction and 100 values of

141    precipitation rates for 5 years, for each $1°$ grid cell that has Terra and Aqua cloud

142    observations.

143

144    **3. Application of *k*-means clustering**

145         In this study we build our basis dataset of hybrid regimes using *k*-means. The *k*-means

146    clustering algorithm (Anderberg 1973; MacQueen 1967) is one of the most popular

147    unsupervised clustering algorithms. This simple algorithm can handle very large data

148    volumes efficiently, hence it is widespread in various studies implementing clustering of

149    geophysical variables, as noted in the Introduction. The underlying principle of the algorithm

7

150   is that for input data consisting of *m_samples* ×*n_features*, feature distances are calculated

151   between each sample and given centroids, and each sample is assigned to the centroid

152   corresponding to the smallest distance. The mean of newly assigned samples becomes the

153   new centroid, and the assignment is repeated until the new centroids are (nearly) identical to

154   the centroids of the previous iteration. Eventually all data are assigned to the group with the

155   most similar members, which minimizes the total Mean Squared Error of the grouped data. In

156   this study, we set the threshold of centroid movement to 1.0e-6, which yields convergence in

157   a few hundred iterations (we set no limit on the total number of iterations).

158   *a. Preparing input data: how to balance between cloud and precipitation data*

159   Previously, Jin et al. (2020) derived tropical cloud regimes (TCRs) using MODIS cloud

160   2D joint histogram data. Since the cloud histogram bin values ranged from 0 to 1 by

161   definition, TCRs could be obtained from the *k*-means clustering algorithm without any

162   normalization process. In order to derive hybrid regimes, the range of values of IMERG

163   precipitation rates must be equivalent to the cloud histogram data. This was easily

164   accomplished by transforming precipitation rates to normalized histogram bin values,

165   similarly to the cloud data.

166   In transforming precipitation data into a histogram, one issue to consider is how to choose

167   the number of bins. Too small a number of bins results in excessive smoothing, which makes

168   notable precipitation patterns indistinguishable. Conversely, too large a number of bins

169   increases noise and prevents us from obtaining meaningful clusters. Since it is known that

170   similar clouds can have varying precipitation rates (e.g., Jin et al. 2018, 2020), we gravitated

171   towards a rather coarser binning of the precipitation histograms. After some testing, we

172   settled on an approximately logarithmically-spaced 6-bin precipitation histogram with bin

173   boundaries at 0.03, 0.1, 0.33, 1, 3.33, 10, 999mm/h. We note that these histogram bin

8

174    boundaries exclude no-rain counts for consistency with the cloud histogram, and also very

175    small precipitation rates below 0.03mm/h.

176        The second issue we had to address was the relative weight between cloud and

177    precipitation when applying the clustering algorithm. If we combine cloud and precipitation

178    histograms without any weighted treatment, the relative importance of cloud compared to

179    precipitation in the *k*-means clustering calculation is 7 to 1 because the cloud histogram

180    consists of 42 bins while the precipitation histogram consists of 6 bins (for a total of 48 bins).

181    With Euclidean distance adopted as the measure to assign data to one of centroids in the *k*-

182    means algorithm, the number of bins translates linearly to relative importance. In this sense, it

183    is possible to make both cloud and precipitation equally important by combining the 42-bin

184    cloud histogram with the precipitation histogram replicated seven times for a total of 84 bins

185    that come from two equal 42-bin contributions from the cloud and precipitation side. In this

186    study, a total of 3 different versions of weights for cloud and precipitation were tested,

187    namely 7:1, 7:3, 7:7.  Only the 7:1 and 7:7 versions will be shown in the manuscript itself,

188    with the 7:3 version shown in the Supplementary Material Part A. We also derive a new set

189    of cloud-only regimes to be used as a reference by following the same procedures, described

190    in the next subsection, as for the hybrid regimes.

191        In terms of regional coverage, we performed the *k*-means algorithm separately for the

192    deep tropics (15°S-15°N) and for much larger portion of the globe that expands to

193    midlatitudes (50°S-50°N). The two domains for 5 years for both Terra and Aqua data result

194    in populations of ~34 million and ~116 million data points once missing values are excluded.

195    In this study, we focus on the deep tropical results only, while the near-global results are

196    shown in the Supplementary Material Part B.

197    *b. Initializing with k-means++ algorithm*

9

198    The *k*-means clustering algorithm is, by definition, deterministic to the initial values,

199    namely the centroids chosen initially. If more than one of the initial centroids are chosen from

200    potentially the same cluster group (i.e., they are similar to each other), the end result of the

201    clustering may not be optimal. To reduce the probability of this happening, and to improve

202    the performance of the *k*-means clustering, a "*k*-means++" algorithm was developed for

203    smarter initialization (Arthur and Vassilvitskii 2007). The *k*-means++ employs a weighted

204    random selection method, where the distance from a pre-selected initial centroid is set as the

205    weight of the data member. If two or more initial centroids are already selected, the minimum

206    distance is selected as the weight. This process ensures that the farthest (largest Euclidean

207    distance) data member from pre-selected centroid(s) has the highest possibility to be chosen,

208    thus ultimately making the initial centroids well-separated from each other. We employ the *k*-

209    means++ algorithm to initialize the *k*-means clustering scheme with 50 different sets of initial

210    centroids (i.e., 50 realizations) for each candidate number *k* of clusters, in order to potentially

211    achieve the best *k*-means clustering results (see next subsection).

212    *c. Criteria for choosing the optimal number of clusters*

213    The *k*-means clustering algorithm requires the number of clusters, *k*, as a preset to be

214    decided by the user. By the nature of *k*-means clustering, a larger number of clusters always

215    decreases the magnitude of "error", measured by the "Within-Cluster(intra-cluster) Variance

216    (WCV)", since the larger *k* the less diverse the members of a group are. At the same time, a

217    large *k* has the undesirable effect of diminishing the level of data compression, which is

218    another way of saying that too many clusters make the grouping less practical and useful. An

219    appropriate value of *k* therefore represents a compromise between the amount of error and the

220    level of compression.

221     Several methods exist to determine the optimal value of $k$. One of the most basic and

222     intuitive methods is the so-called "elbow" method. By observing the percentage of explained

223     variance as a function of the number of clusters, the value of $k$ is selected when the marginal

224     gain of explained variance is small with another cluster added. An issue with this method is

225     that characterizing the gain as marginal is subjective and ambiguous. In many cases the

226     "elbow" point is not obvious, which makes this method unreliable (e.g., Ketchen and Shook

227     1996).

228     The Calinski-Harabasz criterion (CHC; Caliński and Harabasz 1974) is another popular

229     method to determine the most optimal $k$. The basic idea of CHC is to maximize the overall

230     "between-cluster(inter-cluster) variance (BCV)," which indicates maximum separation

231     among clusters, while minimizing the error expressed by WCV. A CHC metric is defined as

232
$$CHC_k = \frac{BCV}{(k-1)} \Big/ \frac{WCV}{(N-k)}$$

233     where $N$ is the total number of data points, and $k$ is the number of clusters. The BCV and

234     WCV are defined as

235
$$BCV = \sum_{i=1}^{k} n_i \|\mu_i - \mu\|^2$$

236
$$WCV = \sum_{i=1}^{k} \sum_{x \in C_i} \|x - \mu_i\|^2$$

237     where $n_i$ is the number of data points in cluster $i$ ($C_i$), $\mu_i$ is the mean of data points in

238     cluster $i$ (a.k.a. centroid), and $\mu$ is the overall mean of all data points. The value of $k$ yielding

239     the maximum CHC represents the best choice for cluster number $k$.

240     The Davies-Bouldin criterion (DBC; Davies and Bouldin 1979) also pursues the

241     maximum separation of clusters with minimum errors in the clusters as CHC, but uses

242     different measures. The DBC metric is defined as

11

File generated with AMS Word template 1.0

$$243 \qquad DBC_k = \frac{1}{k} \sum_{i=1}^{k} max_{j \neq i}\{R_{i,j}\}$$

$$244 \qquad R_{i,j} = \frac{S_i + S_j}{D_{i,j}}$$

245     where $R_{i,j}$ is the ratio of within-cluster scatter of the $i^{th}$ and $j^{th}$ clusters ($S_i$, $S_j$) to the separation

246     between the $i^{th}$ and $j^{th}$ clusters ($D_{i,j}$). $S_i$ and $D_{i,j}$ are defined as

$$247 \qquad S_i = \left( \frac{1}{n_i} \sum_{x \in C_i} \|x - \mu_i\|^2 \right)^{1/2}$$

$$248 \qquad D_{i,j} = \|\mu_i - \mu_j\|$$

249     Here, the within-cluster scatter ($S_i$) represents average distance between each data point and

250     centroid, and the separation measure $D_{i,j}$ is the Euclidean distance between two centroids. For

251     a given $k$, by choosing the maximum ratio for each cluster, DBC measures the worst-case

252     scenario for each cluster. The minimum value of DBC represents therefore the most optimal

253     number of clusters.

254         Figure 1 shows the dependence on $k$ of these criteria in the case of 6 precipitation

255     histogram bins with weight number 1 (i.e., 48-element combined array, referred to as

256     "Cld42+Pr6x1"). The left panel (Fig. 1a) shows maximum BCV and minimum WCV as a

257     function of $k$. The elbow method can be applied to both BCV and WCV. (We note that,

258     because the explained variance is defined as BCV divided by total variance and total variance

259     is a fixed number, it is essentially the same to apply the elbow method to either explained

260     variance or BCV.) However, both BCV and WCV change smoothly as $k$ increases, and it is

261     hard to find an "elbow" in the above figure. In the right panel (Fig. 1b), DBC clearly

262     indicates that 16 is the optimal $k$ while CHC monotonically decreases as $k$ increases. The

263     CHC metric heavily depends on the total population of data points ($N$) by definition, and in

264     the case of huge $N$ ($N \approx$ 34M for our deep tropics domain), variability of CHC is dominated

12

265      by the term, $N/(k-1)$, which results in monotonical decrease with $k$ in a reasonable range.

266      Taking all these into account, DBC is chosen as the primary criterion for selecting the

267      optimal number of clusters, and the trial producing the globally minimum DBC value

268      determines the final set of regimes composed of $k$ centroids. Table 1 shows the values of $k$

269      that came out of this procedure for the four (i.e., including zero) precipitation weights, for

270      both the narrow and extended domains in latitudes. Figures similar to Fig. 1 for the other

271      cases are shown in the Supplementary Materials.

272

273      **4. Details of tropical hybrid regimes**

274      *a. Cloud-only regimes*

275      A set of cloud-only regime is derived as the baseline with which the cloud-precipitation

276      hybrid regimes can be compared to. Jin et al. (2020) previously derived a set of cloud regimes

277      with $k=10$ in the same deep tropics domain, with the last regime being decomposed to 4 sub

278      regimes, for a total of 13 regimes, using the concept of "nested clustering" (Luo et al. 2017;

279      Mason et al. 2014; Oreopoulos et al. 2016). Here, the data period is shortened from 14 years

280      to 5 years to accommodate the availability of precipitation observations, and DBC is

281      employed to select the final set of regimes without invoking nested clustering.

282      Figure 2 shows that the (deep) tropics cloud-only regime (TCR; please note that for

283      economy we drop the tropical "T" designation in the following figures) set is composed of 8

284      high-cloud regimes, 5 low-cloud regimes, and one mixed semi-clear regime (TCR14). Each

285      TCR, except TCR14, has a unique distinct peak of bin cloud fraction. This is a notable

286      difference from the previous TCR set reported by Jin et al. (2020), particularly for high

287      clouds with relatively large optical thickness. Figure 1 in Jin et al. (2020) showed three TCRs

13

288    relevant to convective activity, with peaks of similarly large cloud fraction values at two

289    neighboring histogram bins. These blunt peaks seem to have now split into two TCRs. For

290    example, the old TCR1 had the largest cloud fraction bin values across the cirrostratus (Cs)

291    and cumulonimbus (Cb) bins, according to the traditional ISCCP cloud types (Rossow and

292    Schiffer 1999), and these have now split into peaks that occur in TCR1 and TCR2. By

293    comparing the assignments of each grid cell to old and new TCRs, we confirm that the most

294    grid cells previously assigned to old TCR1, TCR2, and TCR3 in Jin et al. (2020) are now

295    assigned to TCR1 to TCR6. Among them, the first three TCRs dominate precipitation, and

296    TCR1 having the optically thickest and highest cloud dwarfs the other regimes in mean

297    precipitation rate.

298    *b. Hybrid regime2 with precipitation weight of 1 (Cld42+Pr6x1)*

299    We first introduce the tropical cloud-precipitation (hybrid) regime (TCPR) set that

300    corresponds to the precipitation weight of 1 (i.e., cloud-to-precipitation weight ratio is 7:1

301    with 48-element array; Cld42+Pr6x1). By adding precipitation information this way, the

302    optimal number of clusters according to the DBC increases from 14 to 16 in our tropical

303    domain (Table 1 and Fig. 3). This TCPR set is composed of 9 high cloud regimes, 5 low

304    cloud regimes, and 2 mixed regimes (including a semi-clear regime, TCPR16). A notable

305    difference in centroids when rainfall information added is the newly occurring TCPR10. This

306    regime represents high and low mixed clouds with intermediate cloud fraction and substantial

307    precipitation. In order to investigate the origin of this version of TCPR10, we introduce a

308    regime coincidence distribution matrix (Fig. 4) showing the RFO of new regimes (i.e.,

309    Cld42+Pr6x1; x-axis) for the grid cells assigned to one of the cloud-only regimes (y-axis).

310    This graphical matrix indicates that grid cells assigned to TCPR10 belonged previously to

311    various TCRs (e.g., TCR3, 5, 7, 10, 14, etc.). In terms of population, the biggest contributor

14

312    is TCR14 which is semi-clear regime with RFO 38% (Fig. 2). Because of the split of TCR14

313    due to the addition of rainfall information, the similar semi-clear hybrid regime (TCPR16)

314    has now a lower RFO value (32.8% in Pr6x1 TCPR16 vs. 37.9% for TCR14) and a lower

315    cloud fraction (26% vs. 32%).

316        The other contributor to the increased $k$ from the cloud-only regimes to hybrid regimes is

317    the split of TCR8 into TCPR8 and TCPR9. TCR8 in Fig. 2 represents a cirrus (Ci)-dominant

318    regime with a cloud fraction peak in the bin of highest cloud top (lowest CTP) and smallest

319    optical thickness; it is now split into two versions of Ci-dominated regimes with total cloud

320    fractions of 58% (TCPR8) and 78% (TCPR9). While neither TCPR8 nor TCPR9 seem to be

321    producing substantial rainfall, the precipitation histogram component of the centroid shows

322    that TCPR8 has a slightly elevated chance of intermediate intensity precipitation.

323        It is also worth noting that significant fractions of grid cells occupied by TCR3 are now

324    assigned to TCPR5 in addition to TCPR3. TCPR3 and TCPR5 show clearly different

325    precipitation characteristics: the estimated average precipitation rate of TCPR3 is 1.2mm/h

326    with the peak of precipitation histogram around 1mm/h while the average rate of TCPR5 is

327    0.2mm/h. A possible interpretation is that TCR3 has grid cells of similar clouds with varying

328    precipitation intensities from light to intermediate, and grid cells of lighter precipitation are

329    shifted to TCPR5 by the addition of precipitation information. Similar phenomena of lighter

330    rain grid cells shifted to other hybrid regimes are also found for TCR1, TCR5, and TCR6

331    indicating that within-regime precipitation variability decreases in the hybrid regimes because

332    outliers with weak precipitation in cloud-only regimes are now removed. On the other hand,

333    regimes dominated by low clouds show great consistency between the cloud-only and hybrid

334    regime sets because there are barely any precipitation features that would make them

335    distinguishable.

15

*c. Hybrid regimes with precipitation weight of 7 (Cld42+Pr6x7)*

337      As the relative weight of precipitation increases from 1 to 3, the patterns of the cloud joint

338    histogram component of the centroids lose peak sharpness, and some regimes even show

339    blunt peaks across two adjacent levels of CTP (see Supplementary Material Part A). As the

340    relative weight of precipitation further increases to 7, namely when cloud and precipitation

341    histograms matter equally in the (84-element) combined arrays subjected to $k$-means

342    clustering, the patterns of the mean joint cloud histogram exhibit even blunter peaks, and

343    some hybrid regimes now share quite similar cloud patterns (e.g., TCPR3 and TCPR4;

344    TCPR7 and TCPR8 and TCPR9 in Fig. 5). This suggests that precipitation rather than cloud

345    has now a greater impact in determining the assignment to certain TCPRs, and a previous

346    regime of the no or small precipitation weight set can be split into multiple regimes

347    depending on the shape of the precipitation histogram. Indeed, the optimal number of clusters

348    in the Cld42+Pr6x7 case (a.k.a. equal-weight set) increases to 19, with 13 high cloud

349    regimes, 4 low cloud regimes and 2 mixed regimes (including the semi-clear regime).

350      A similar matrix of regime coincidence distribution between precipitation weight number

351    1 and 7 is displayed in Fig. 6. In the Cld42+Pr6x1 set, TCPR1, TCPR2 and TCPR3

352    represents high and thick clouds producing intermediate to heavy precipitation. All these

353    three TCPRs are now split into 3 or more TCPRs in the equal-weight set because of the

354    increased impact of precipitation on the clustering. As a result, centroids of equal-weight set

355    show distinct patterns in the precipitation histogram component of the centroid, something

356    that can be interpreted as decreased variability in precipitation intensity and increased

357    variability in cloud type mixtures in the grid cells belonging to a specific TCPR of the equal-

358    weight set. Also noteworthy is that TCPR10 of Cld42+Pr6x1 which was diagnosed as

359    representing mixed clouds with intermediate precipitations is now split into 4 different

360    TCPRs. In terms of cloud histogram pattern, TCPR14 of Fig. 5 shares some similarity with

361    TCPR10 of Fig. 3, but TCPR14 of the equal-weight set has notably smaller high-cloud

362    fractions and intermediate-precipitation fractions. The decomposition of TCPR10 in

363    Cld42+Pr6x1 is a major contributor to the increased number of clusters from 16 to 19.

364         In summary, we find that the information added by precipitation helps to also distinguish

365    clouds with a greater degree of detail in terms of cloud-precipitation relationship. In the set

366    where the added precipitation information matters the least, namely the 7:1 weight ratio

367    (Cld42+Pr6x1), the cloud histogram patterns are mostly consistent with the cloud-only

368    regimes. Still, the added precipitation information rearranges some outlier grid cells in cloud-

369    only regimes (in terms of precipitation properties), thus resulting in tighter relationships

370    between cloud and precipitation in the new regimes. The enhanced weight of precipitation

371    obviously decreases the influence of cloud patterns in the resulting centroids, even to the

372    degree where similar cloud histogram patterns (albeit with distinct precipitation histogram

373    patterns) appear in the equal-weight set. These cloud and precipitation pattern changes occur

374    mostly in regimes dominated by high-clouds; regimes dominated by low clouds are not

375    changing much by increasing the precipitation weight indicating the lack of diversity in

376    precipitation properties, at least according to IMERG.

377

378    **5. Projection onto IMERG domain**

379    *a. Can cloud be predicted from precipitation?*

380         Cloud and precipitation are closely related, but at the same time there is significant

381    precipitation variability within similar clouds, and vice versa. In the previous section, we

382    showed two sets of tropical cloud-precipitation hybrid regimes, representing the dominant

17

383    mixtures of specific cloud types and corresponding precipitation intensities (other variants of

384    relative weights and an extension that includes extratropics are shown in the Supplemental

385    Materials). In this section, we examine the feasibility of "predicting" clouds from solely

386    precipitation information using these hybrid regimes. The reason we want to predict clouds is

387    because the cloud observations suffer from substantial amounts of missing grid cells due to

388    the swath width of the MODIS granules, and are much sparser temporally compared to the

389    IMERG precipitation dataset. An extended dataset of cloud information with higher temporal

390    resolution could be useful for various research endeavors.

391        The availability of cloud-precipitation hybrid regimes simplifies a potential cloud

392    prediction scheme because clouds in a grid cell are represented by the limited number of

393    classes (regimes) derived from the clustering analysis. (Additional information about the

394    clouds besides what hybrid regime they belong would obviously not be available.) Hence, the

395    problem at hand is predicting one of the hybrid regimes based on only the precipitation

396    information of a grid cell. The simplest way to assign a hybrid regime to grid cell at a time

397    when no cloud information is available is to adopt the Euclidean distance criterion used in the

398    $k$-mean clustering, but now applied only on the observed IMERG precipitation histogram and

399    the precipitation component of the hybrid regime centroid. Of course, this assignment by

400    precipitation is only possible when a reasonable amount of precipitation is detected;

401    identification of hybrid regime occurrence in a grid cell where barely any rain occurs is

402    impossible.

403        The performance of hybrid regime prediction by matching observed and centroid

404    precipitation histograms is summarized in Fig. 7 for the case of the equal-weight set

405    (Cld42+Pr6x7) in the extended tropical domain of 20°S to 20°N. Figure 7 is a Fig.4-like

406    regime coincidence distribution matrix between original TCPRs (y-axis) observed at the time

File generated with AMS Word template 1.0

407    of Terra and Aqua daytime overpasses and predicted TCPRs by precipitation-only (x-axis)

408    for the same grid cells. Among the 19 regimes, those with precipitation fraction (= sum of 6

409    bins of precipitation histogram) below 10% are merged into the "Others" class. Overall, the

410    prediction results are quite impressive; among regimes having significant amounts of

411    precipitation, five regimes (CPR1, 2, 4, 8, 9) have precipitation-based prediction accuracy of

412    TCPR occurrence above 95%. Furthermore, the accuracy of TCPR3 and TCPR7 prediction is

413    also quite high, over 90%. This means that the precipitation signatures of members belonging

414    to these hybrid regimes are unique enough to allow them to be differentiated from members

415    of other regimes. These regimes commonly have precipitation fractions above 50%. While

416    for TCPR5, exhibiting only 20% prediction accuracy, the estimated mean precipitation is

417    greater than that of TCPR9, precipitation fraction is just 23%, less than half of TCPR9's (Fig.

418    5). A small total precipitation fraction usually means that histogram bin values are also small,

419    which makes them hard to be distinguished from other regimes under our adopted Euclidean

420    distance criterion. In addition, we also examined the accuracies geographically (by

421    longitudes), and found that prediction accuracies are quite stable regardless of longitudes

422    with only small drops of accuracy in the central Africa and South America in the case of

423    TCPR1 and TCPR7 (see Supplementary Material Part A).

424        The equal-weight set shows that the regimes having intermediate-to-heavy precipitation

425    intensity can by predicted well by the precipitation-only histogram constructed by the 0.1°

426    IMERG data, a result likely due to the significant impact of precipitation on the clustering

427    process. We also tested the case of small precipitation weight, and as expected, the prediction

428    accuracy was markedly lower, as shown in Fig. 8. In the case of Cld42+Pr6x1, 7 regimes

429    pass the criterion of precipitation fraction above 10% among the 16 regimes. The highest

430    accuracy, 81%, is achieved by TCPR1 which has the heaviest precipitation and thickest

431  clouds representing a group of convective cores (Fig. 3). The second highest prediction

432  accuracy, 50%, is achieved by TCPR10 (mixed cloud types and a large fraction of light rain),

433  while all other prediction skills are below 50%. In the case of TCPR2 and TCPR3 both of

434  which have intermediate precipitation intensity, precipitation histogram patterns are too

435  similar (Fig. 3) for them to be separable in the regime prediction. Still, Fig. 8 suggests a hint

436  of different precipitation characteristics between TCPR2 and TCPR3, where the light

437  precipitations tails of TCPR2's rainfall distribution gives rise to TCPR10 assignment for 14%

438  of the grid cells, while TCPR3 being biased towards heavy precipitation results in assignment

439  of 13% of the grid cells to TCPR1.

440      To summarize, we demonstrated that we can predict cloud patterns through the prediction

441  of hybrid regimes from precipitation-only information when using the set of hybrid regimes

442  derived with equal weighting between cloud and precipitation (Cld42+Pr6x7). A total of 7

443  hybrid regimes can be predicted highly accurately when their precipitation features include

444  intermediate to heavy rainfall intensity and their cloudiness corresponds to high-thick cloud

445  patterns. In practical terms this means that through the process of assigning regimes by

446  precipitation histogram Euclidean distance, we can transform the 30-minute full tropical

447  coverage IMERG data into occurrence maps of these 7 regimes at 1-degree resolution and at

448  the same 30-minute temporal resolution, i.e. we have achieved a projection of TCPRs onto

449  the IMERG domain. In the following subsection, we present an application example of this

450  newly built hybrid regime occurrence maps.

451  *b. Analysis example: Diurnal cycle of hybrid regimes*

452      Due to the reliance of cloud optical thickness retrievals on the availability of solar

453  insolation, 2D joint histogram data of cloud is available once daily for each of Terra and

454  Aqua, at around 10:30am and 1:30pm local solar time (LST), respectively. Hence, even a

20

455  combined analysis of Terra and Aqua can provide only limited information on cloud

456  variability around noon in LST. The occurrence map of hybrid regimes projected onto

457  IMERG domain according to our method described in the previous subsection radically

458  improves the temporal resolution (30-min), thus enabling examination of the diurnal cycle of

459  the hybrid regimes for which we have good prediction capability, based on the assumption

460  that nighttime cloud-precipitation relationship remains the same as in daytime. Figures 9 and

461  10 show the RFO of TCPR1 and TCPR2 of the Cld42+Pr6x7 set in the longitude-LST phase

462  space, respectively. We note that LST is calculated by adding the regionally-dependent

463  factor, longitude×(24/360) to UTC as in Tan et al. (2019b).

464      TCPR1 of the Cld42+Pr6x7 set represents deep convective cores with the heaviest

465  precipitation. Previously, Jin et al. (2018, 2020) showed that the regime corresponding to the

466  heaviest precipitation most frequently occurs in the tropical warm pool oceans. Figure 9 is

467  consistent with the previous studies, and shows the highest RFO in the east and west of the

468  Maritime Continent. Moreover, the temporal evolution indicates that the most active hour of

469  TCPR1 occurrence is in the early morning, 2am to 8am in this region, consistent with Fig. 11

470  in Yang and Smith (2006), but deviating from the findings of Kikuchi and Wang (2008) who

471  noted oceanic peak between 6am to 9am. Other than the warm pool region, TCPR1 also

472  notably occurs in the Amazon basin, and is slightly more active in the early morning than

473  other local times, which is consistent with the precipitation diurnal cycle driven by dynamical

474  processes (Vernekar et al. 2003). Regardless of the longitude, a hint of local RFO minimum

475  appears just before noon, a feature that actually becomes clearer when examining TCPR2 in

476  Fig. 10.

477      TCPR2 of the Cld42+Pr6x7 set also responds to quite heavy precipitation, with the peak

478  of cloud fraction occurring at the same CTP level, but for slightly optically thinner clouds

479    (Fig. 5), suggesting a combination of convective cores and thick anvils. Figure 10 shows that

480    TCPR2 also frequently occurs in the tropical warm pool oceans and Amazon basin, like

481    TCPR1. However, the active hours are clearly different from TCPR1. For example, in

482    addition to the early morning times as in Fig. 9, TCPR2 also frequently occurs just before

483    noon and in the afternoon between 2pm and 6pm in the warm pool region. In the Amazon

484    basin, the most active hour is shifted to afternoon, the time previous studies noted as the most

485    active hours of continental convection driven by thermodynamic processes (Giles et al. 2020;

486    Janowiak et al. 2005).

487       In Fig. 10 we can see RFO local minima troughs four times a day: 12am-2am, 8am-10am,

488    12pm-2pm, and 8pm-10pm. In these time windows, the occurrence of TCPR2 decreases

489    abruptly, which may suggest an artifact in the IMERG dataset. Similar trough-like patterns

490    are also detected with other TCPRs, notably for TCPR4, TCPR8, and TCPR9, as for TCPR2,

491    and less prominently for TCPR3 and TCPR7, as for TCPR1 (see Supplementary Material Part

492    A). The troughs, especially spaced in two pairs 12-h apart, points to the possibility of an

493    artifact stemming from particular sensors on board sun-synchronous satellites used in

494    IMERG. In particular, these times match the overpass times of several cross-track scanning

495    sounders in the constellation which generate double-peaks in precipitation distributions over

496    ocean (You et al. 2020). However, troughs of the same diurnal cycle analysis over land-only

497    are still notable (but with weakened signal; see Supplementary Material Part A), indicating

498    that there may be other unidentified factors at play or that the troughs represent true diurnal

499    signals in fact.

500       In summary, through the projection of hybrid regimes onto IMERG domain, the temporal

501    resolution for some of regimes with the greatest precipitation contribution and most likely

502    associated with convection, is greatly improved. In addition to the diurnal cycle analysis

503 shown in this subsection, this diurnally-extended dataset of cloud-precipitation hybrid

504 regimes has enormous potential to examine other features of convective systems. We should

505 also note that this projection method works not only for the deep tropical regimes, but also

506 for the hybrid regimes of extended latitudes when the higher weights of precipitation are used

507 in the clustering procedure (see Supplementary Material Part B).

508

## 6. Summary and Conclusion

510     We generated hybrid cloud and precipitation regimes (CPRs) by applying the $k$-means

511 clustering algorithm, with advanced initialization and objective measures to select the optimal

512 number of clusters $k$, on coincident cloud and precipitation data from MODIS and IMERG.

513 We discussed how multiple versions of hybrid CPR sets can be obtained depending on the

514 relative weighting of the cloud and precipitation information and the boundaries of the

515 geographical domain.

516     Given that precipitation was represented by a rather coarse 6-bin histogram and clouds

517 were represented by a 42-element joint histogram, a naïve concatenation of cloud and

518 precipitation arrays implies a 7 to 1 ratio in cloud versus precipitation weighting. When

519 performing joint clustering with this 48-element array, the patterns of the cloud histogram

520 centroids looked quite similar to those of cloud-only centroids, indicating a weak influence of

521 precipitation on the clustering. However, for the cloud regime associated with intermediate to

522 heavy rainfall intensity, some outliers with relatively lighter rainfall were moved to other

523 regimes of corresponding rainfall intensity, making the precipitation variability of hybrid

524 regimes generally tighter. As the weight of precipitation in the joint clustering progressively

525 increased (by replicating the precipitation histograms as needed), the precipitation histogram

23

526 component of the hybrid centroids became more unique from those of the other centroids

527 while the cloud histogram parts of the centroids started losing peak sharpness. In the set of

528 equal-weight between cloud and precipitation, three CPRs of high clouds with light-to-

529 intermediate precipitation intensity even shared quite similar cloud histogram patterns (but

530 with distinct precipitation histogram patters, of course). Compared to the high cloud regimes

531 experiencing dramatic changes by varying the weight of precipitation, low cloud regimes

532 remained relatively unchanged among different sets, because their weak rainfall did not

533 impact the clustering process.

534     Given that the precipitation histogram part of centroid became progressively more distinct

535 from that of the other centroids as precipitation weight increased, we tested whether we can

536 predict a specific CPR based on only the precipitation information of the grid cell. This

537 attempt was motivated by the fact that IMERG dataset has much higher temporal resolution

538 with nearly no missing data at 30-minute intervals compared to temporally sparse MODIS

539 cloud observations. We found that, in the case of equal-weight set, seven high cloud regimes

540 with intermediate-to-heavy precipitation can be predicted with over 90% accuracy by the

541 precipitation information only. This result suggests that a projection of certain CPRs onto the

542 IMERG domain is possible, opening thus a broad path for a variety of studies that require

543 diurnally-resolved cloud information.

544     In a previous study by Jin et al. (2020), three cloud-only regimes related to tropical

545 convective activities were selected, to study various features of convective systems at

546 synoptic scales. However, their investigation was limited to snapshots of convective systems

547 near 1:30pm LST due to the limitation of MODIS cloud observation availability and with

548 morning Terra observations filling swath gaps based on persistence assumptions. The

549 IMERG-based projection method enabled by hybrid regimes as mentioned above can expand

550     their study in various directions. For example, thanks to 30-minute temporal resolution

551     without gaps, diurnal cycle of convective systems can be examined in a manner demonstrated

552     in Figs. 9 and 10. In addition, it is also possible to examine the life cycle of large-scale

553     convective systems by systematically tracking them. While the prediction skill using IMERG

554     precipitation is not perfect at all instances, the expansion of hybrid regimes to temporally

555     high resolution is a significant advancement that can contribute to better understanding of

556     large-scale tropical convective systems.

557

563

564     *Data Availability Statement.*

565       *IMERG precipitation data used in this study is openly available from the NASA Goddard*

566     *Earth Sciences Data and Information Services Center (GES DISC) at*

567     https://doi.org/10.5067/GPM/IMERG/3B-HH/06 as cited in Huffman et al. (2019b). *Daily*

568     *MODIS L3 cloud histogram data for Terra (MOD08_D3) and Aqua (MYD08_D3) are openly*

569     *available from the Level-1 and Atmosphere Archive & Distribution System (LAADS)*

570     *Distributed Active Archive Center (DAAC) in the Goddard Space Flight Center at*

571     *https://doi.org/10.5067/MODIS/MOD08_D3.061 and*

572     *https://doi.org/10.5067/MODIS/MYD08_D3.061 as cited in Platnick et al. (2017b). The*

573     *MODIS cloud regime and MODIS-IMERG cloud-precipitation hybrid regime datasets*

25

File generated with AMS Word template 1.0

574 *derived in 15°S-15°N domain is available at https://data.nasa.gov/Earth-Science/Cloud-*

575 *Precipitation-Hybrid-Regimes-MODIS-IMERG-in-/ee3g-swmf.*

576

577                                          REFERENCES

578 Anderberg, M. R., 1973: *Cluster Analysis for Applications*. Elsevier, 359 pp.

579 Arthur, D., and S. Vassilvitskii, 2007: k-means++: the advantages of careful seeding. *SODA 07*

580     *Proc. Eighteenth Annu. ACM-SIAM Symp. Discrete Algorithms*, 1027–1035.

581 Caliński, T., and J. Harabasz, 1974: A dendrite method for cluster analysis. *Commun. Stat. -*

582     *Theory Methods*, **3**, 1–27, https://doi.org/10.1080/03610927408827101.

583 Davies, D. L., and D. W. Bouldin, 1979: A Cluster Separation Measure. *IEEE Trans. Pattern*

584     *Anal. Mach. Intell.*, **PAMI-1**, 224–227, https://doi.org/10.1109/TPAMI.1979.4766909.

585 Giles, J. A., R. C. Ruscica, and C. G. Menéndez, 2020: The diurnal cycle of precipitation over

586     South America represented by five gridded datasets. *Int. J. Climatol.*, **40**, 668–686,

587     https://doi.org/10.1002/joc.6229.

588 Huffman, G. J., and Coauthors, 2019a: Algorithm Theoretical Basis Document (ATBD) version

589     06. NASA Global Precipitation Measurement (GPM) Integrated Multi-satellitE Retrievals

590     for GPM (IMERG).

591     https://gpm.nasa.gov/sites/default/files/document_files/IMERG_ATBD_V06_0.pdf.

592 ——, E. F. Stocker, D. T. Bolvin, E. J. Nelkin, and J. Tan, 2019b: GPM IMERG Final

593     Precipitation L3 Half Hourly 0.1 degree x 0.1 degree V06, Greenbelt, MD, Goddard Earth

594     Sciences Data and Information Services Center (GES DISC).

595     https://doi.org/10.5067/GPM/IMERG/3B-HH/06.

File generated with AMS Word template 1.0

596     Jakob, C., and G. Tselioudis, 2003: Objective identification of cloud regimes in the Tropical

597         Western Pacific. *Geophys. Res. Lett.*, **30**, 2082, https://doi.org/10.1029/2003GL018367.

598     Janowiak, J. E., V. E. Kousky, and R. J. Joyce, 2005: Diurnal cycle of precipitation determined

599         from the CMORPH high spatial and temporal resolution global precipitation analyses. *J.*

600         *Geophys. Res.*, **110**, D23105, https://doi.org/10.1029/2005JD006156.

601     Jin, D., L. Oreopoulos, D. Lee, N. Cho, and J. Tan, 2018: Contrasting the co-variability of

602         daytime cloud and precipitation over tropical land and ocean. *Atmospheric Chem. Phys.*,

603         **18**, 3065–3082, https://doi.org/10.5194/acp-18-3065-2018.

604     ——, ——, ——, J. Tan, and K. Kim, 2020: Large-Scale Characteristics of Tropical Convective

605         Systems Through the Prism of Cloud Regime. *J. Geophys. Res. Atmospheres*, **125**,

606         e2019JD021157, https://doi.org/10.1029/2019JD031157.

607     Ketchen, D. J., and C. L. Shook, 1996: The Application of Cluster Analysis in Strategic

608         Management Research: An Analysis And Critique. *Strateg. Manag. J.*, **17**, 441–458,

609         https://doi.org/10.1002/(SICI)1097-0266(199606)17:6<441::AID-SMJ819>3.0.CO;2-G.

610     Kikuchi, K., and B. Wang, 2008: Diurnal Precipitation Regimes in the Global Tropics. *J. Clim.*,

611         **21**, 2680–2696, https://doi.org/10.1175/2007JCLI2051.1.

612     King, M. D., and Coauthors, 2003: Cloud and aerosol properties, precipitable water, and

613         profiles of temperature and water vapor from MODIS. *IEEE Trans. Geosci. Remote Sens.*,

614         **41**, 442–458, https://doi.org/10.1109/TGRS.2002.808226.

615     Lee, D., L. Oreopoulos, G. J. Huffman, W. B. Rossow, and I.-S. Kang, 2013: The Precipitation

616         Characteristics of ISCCP Tropical Weather States. *J. Clim.*, **26**, 772–788,

617         https://doi.org/10.1175/JCLI-D-11-00718.1.

618    Luo, Z. J., R. C. Anderson, W. B. Rossow, and H. Takahashi, 2017: Tropical cloud and

619        precipitation regimes as seen from near-simultaneous TRMM, CloudSat, and CALIPSO

620        observations and comparison with ISCCP: Tropical Clouds From Radars and Lidar. *J.*

621        *Geophys. Res. Atmospheres*, **122**, 5988–6003, https://doi.org/10.1002/2017JD026569.

622    MacQueen, J., 1967: Some methods for classification and analysis of multivariate

623        observations. *Proceedings of the fifth Berkeley symposium on mathematical statistics*

624        *and probability*, Vol. 1 of, Oakland, CA, USA., 281–297.

625    Mason, S., C. Jakob, A. Protat, and J. Delanoë, 2014: Characterizing Observed Midtopped

626        Cloud Regimes Associated with Southern Ocean Shortwave Radiation Biases. *J. Clim.*, **27**,

627        6189–6203, https://doi.org/10.1175/JCLI-D-14-00139.1.

628    Oreopoulos, L., and William. B. Rossow, 2011: The cloud radiative effects of International

629        Satellite Cloud Climatology Project weather states. *J. Geophys. Res.*, **116**, D12202,

630        https://doi.org/10.1029/2010JD015472.

631    ——, N. Cho, D. Lee, S. Kato, and G. J. Huffman, 2014: An examination of the nature of

632        global MODIS cloud regimes. *J. Geophys. Res. Atmospheres*, **119**, 8362–8383,

633        https://doi.org/10.1002/2013JD021409.

634    ——, ——, ——, and ——, 2016: Radiative effects of global MODIS cloud regimes. *J.*

635        *Geophys. Res. Atmospheres*, **121**, 2299–2317, https://doi.org/10.1002/2015JD024502.

636    Pike, M., and B. R. Lintner, 2020: Application of Clustering Algorithms to TRMM Precipitation

637        over the Tropical and South Pacific Ocean. *J. Clim.*, **33**, 5767–5785,

638        https://doi.org/10.1175/JCLI-D-19-0537.1.

28

639     Pincus, R., S. Platnick, S. A. Ackerman, R. S. Hemler, and R. J. P. Hofmann, 2012: Reconciling

640         Simulated and Observed Views of Clouds: MODIS, ISCCP, and the Limits of Instrument

641         Simulators. *J. Clim.*, **25**, 4699–4720, https://doi.org/10.1175/JCLI-D-11-00267.1.

642     Platnick, S., M. D. King, S. A. Ackerman, W. P. Menzel, B. A. Baum, J. C. Riedi, and R. A. Frey,

643         2003: The MODIS cloud products: algorithms and examples from terra. *IEEE Trans.*

644         *Geosci. Remote Sens.*, **41**, 459–473, https://doi.org/10.1109/TGRS.2002.808301.

645     ——, and Coauthors, 2017a: The MODIS Cloud Optical and Microphysical Products:

646         Collection 6 Updates and Examples From Terra and Aqua. *IEEE Trans. Geosci. Remote*

647         *Sens.*, **55**, 502–525, https://doi.org/10.1109/TGRS.2016.2610522.

648     ——, M. D. King, and P. A. Hubanks, 2017b: *MODIS Atmosphere L3 Daily Product (C6.1)*.

649         NASA MODIS Adaptive Processing System, Goddard Space Flight Center,

650         [doi:10.5067/MODIS/MOD08_D3.061; doi:10.5067/MODIS/MYD08_D3.061],.

651     ——, and Coauthors, 2018: MODIS cloud optical properties: User guide for the collection

652         6/6.1 Level-2 MOD06/MYD06 product and associated Level-3 datasets, Version 1.1.

653         https://atmosphere-

654         imager.gsfc.nasa.gov/sites/default/files/ModAtmo/MODISCloudOpticalPropertyUserGui

655         deFinal_v1.1_1.pdf.

656     Robertson, A. W., N. Vigaud, J. Yuan, and M. K. Tippett, 2020: Toward Identifying

657         Subseasonal Forecasts of Opportunity Using North American Weather Regimes. *Mon.*

658         *Weather Rev.*, **148**, 1861–1875, https://doi.org/10.1175/MWR-D-19-0285.1.

659    Rossow, W. B., and R. A. Schiffer, 1999: Advances in Understanding Clouds from ISCCP. *Bull.*

660        *Am. Meteorol. Soc.*, **80**, 2261–2287, https://doi.org/10.1175/1520-

661        0477(1999)080<2261:AIUCFI>2.0.CO;2.

662    ——, G. Tselioudis, A. Polak, and C. Jakob, 2005: Tropical climate described as a distribution

663        of weather states indicated by distinct mesoscale cloud property mixtures. *Geophys.*

664        *Res. Lett.*, **32**, L21812, https://doi.org/10.1029/2005GL024584.

665    ——, A. Mekonnen, C. Pearl, and W. Goncalves, 2013: Tropical Precipitation Extremes. *J.*

666        *Clim.*, **26**, 1457–1466, https://doi.org/10.1175/JCLI-D-11-00725.1.

667    Tan, J., and L. Oreopoulos, 2019: Subgrid Precipitation Properties of Mesoscale Atmospheric

668        Systems Represented by MODIS Cloud Regimes. *J. Clim.*, **32**, 1797–1812,

669        https://doi.org/10.1175/JCLI-D-18-0570.1.

670    ——, C. Jakob, W. B. Rossow, and G. Tselioudis, 2015: Increases in tropical rainfall driven by

671        changes in frequency of organized deep convection. *Nature*, **519**, 451–454,

672        https://doi.org/10.1038/nature14339.

673    ——, G. J. Huffman, D. T. Bolvin, and E. J. Nelkin, 2019a: IMERG V06: Changes to the

674        Morphing Algorithm. *J. Atmospheric Ocean. Technol.*, **36**, 2471–2482,

675        https://doi.org/10.1175/JTECH-D-19-0114.1.

676    ——, ——, ——, and ——, 2019b: Diurnal Cycle of IMERG V06 Precipitation. *Geophys. Res.*

677        *Lett.*, **46**, 13584–13592, https://doi.org/10.1029/2019GL085395.

678    Tselioudis, G., W. Rossow, Y. Zhang, and D. Konsta, 2013: Global Weather States and Their

679        Properties from Passive and Active Satellite Cloud Retrievals. *J. Clim.*, **26**, 7734–7746,

680        https://doi.org/10.1175/JCLI-D-13-00024.1.

File generated with AMS Word template 1.0

681    Vernekar, A. D., B. P. Kirtman, and M. J. Fennessy, 2003: Low-Level Jets and Their Effects on

682        the South American Summer Climate as Simulated by the NCEP Eta Model. *J. Clim.*, **16**,

683        297–311, https://doi.org/10.1175/1520-0442(2003)016<0297:LLJATE>2.0.CO;2.

684    Yang, S., and E. A. Smith, 2006: Mechanisms for Diurnal Variability of Global Tropical Rainfall

685        Observed from TRMM. *J. Clim.*, **19**, 5190–5226, https://doi.org/10.1175/JCLI3883.1.

686    You, Y., V. Petkovic, J. Tan, R. Kroodsma, W. Berg, C. Kidd, and C. Peters-Lidard, 2020:

687        Evaluation of V05 Precipitation Estimates from GPM Constellation Radiometers Using

688        KuPR as the Reference. *J. Hydrometeorol.*, **21**, 705–728, https://doi.org/10.1175/JHM-D-

689        19-0144.1.

690    Zhang, Z., and S. Platnick, 2011: An assessment of differences between cloud effective

691        particle radius retrievals for marine water clouds from three MODIS spectral bands. *J.*

692        *Geophys. Res.*, **116**, https://doi.org/10.1029/2011JD016216.

693    Zhao, J., J. Kug, J. Park, and S. An, 2020: Diversity of North Pacific Meridional Mode and Its

694        Distinct Impacts on El Niño-Southern Oscillation. *Geophys. Res. Lett.*, **47**,

695        e2020GL088993, https://doi.org/10.1029/2020GL088993.

696

697                          TABLES

698    **Table 1**. Optimal values of *k* according to the DBC metric for the two domains and four
699    precipitation weights.

| Deep Tropics (15°S-15°N) | | Low-to-Mid Latitudes (50°S-50°N) | |
|---|---|---|---|
| Cloud-only | k=14 | Cloud-only | k=15 |
| Pr_wt=1 | k=16 | Pr_wt=1 | k=20 |

| Pr_wt=3 | k=16 | Pr_wt=3 | k=19 |
|---------|------|---------|------|
| Pr_wt=7 | k=19 | Pr_wt=7 | k=22 |

700

701 <div align="center">FIGURES</div>



702

703 Figure 1. Criteria for selecting optimal number of clusters (*k*) are displayed as a function of *k*
704 for the case of 7:1 weighting in the combined cloud-precipitation array (Cld42+Pr6x1). (a)
705 Between-cluster variance (BCV; blue circles) and within-cluster variance (WCV; orange
706 triangles), and (b) Calinski-Harabasz criterion (CHC; green circles) and Davies-Bouldin
707 criterion (DBC; red triangles). We note that for the same *k*, a set of initial centroids (i.e., one
708 realization) selected as the best by one criterion can be different from that selected for
709 another criterion.

710

32

Figure 2. Deep tropics cloud-only regime centroids (mean histograms, left) and geographical distribution of relative frequency of occurrence (RFO, right). Bin cloud fraction values exceeding 5% are shown explicitly on the centroid panels. The precipitation histograms shown below the cluster centroids are composite means for each cloud regime. In addition to the total cloud fraction, total precipitation fraction which is the sum of all precipitation histogram bin values, and estimated mean precipitation rate based on the histogram are also given on the panel title. Above the RFO panels, individual Terra and Aqua RFOs are provided in brackets.

33

Figure 3. Similar to Fig. 2 but now the precipitation part of the centroids was also derived from clustering where precipitation contributed with a weight number of 1 (Cld42+Pr6x1; i.e., 7:1 ratio in 48-element combined array in clustering).

Figure 4. Regime coincidence distribution matrix comparing assignment frequencies on the same grid cell between the cloud-only regimes of Fig. 2 (y-axis) and the hybrid regimes Cld42+Pr6x1 of Fig. 3 (x-axis). The values of the matrix are normalized across rows, and values above 10% are explicitly shown. Please note that while the regimes were derived with data in 15°S-15°N, regime assignment was performed in the extended domain, 20°S-20°N for both Terra and Aqua, because tropical phenomena often extend beyond the 15° latitude boundaries.

35

Figure 5. As Fig. 3 but with precipitation contributing with weight number 7 (Cld42+Pr6x7; i.e., 7:7 ratio in 84-element combined array in clustering).

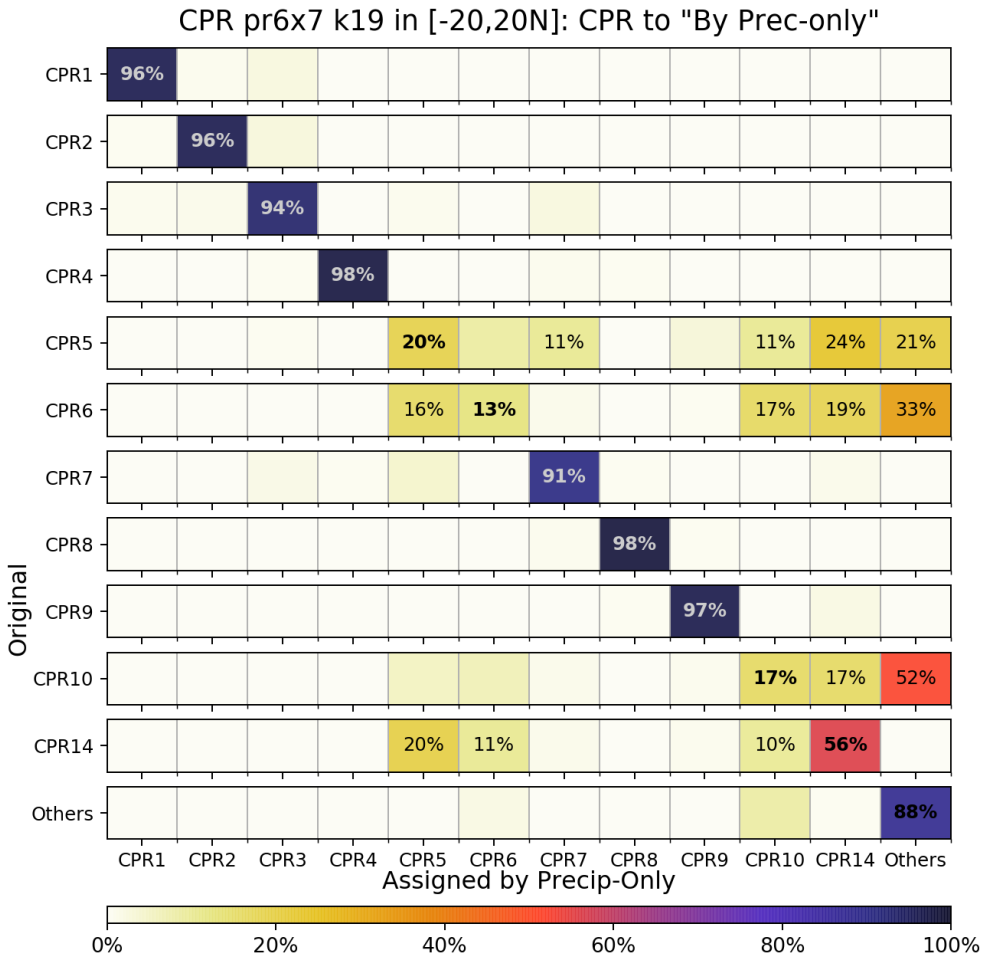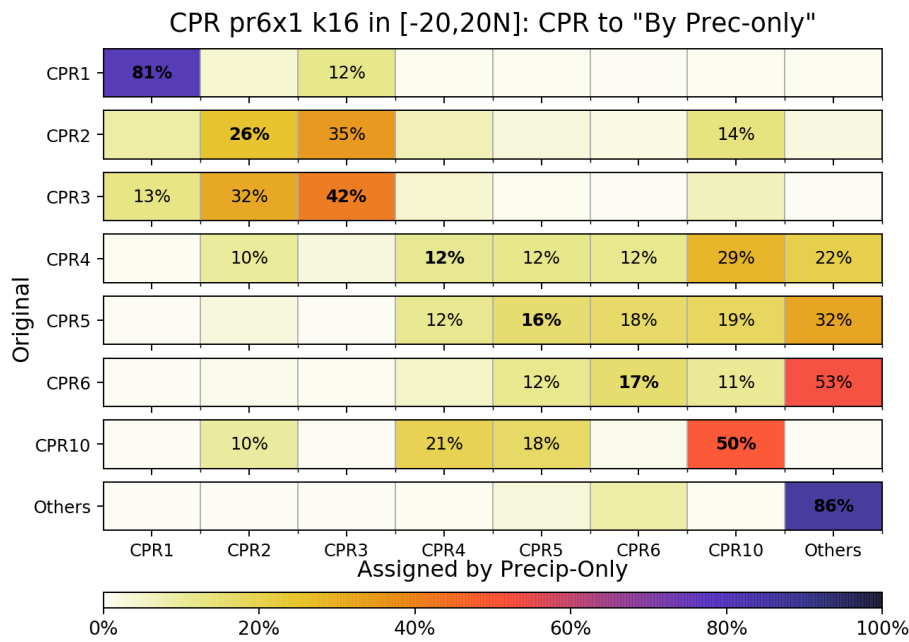Figure 6. As Figure 4 but between Cld42+Pr6x1 (y-axis; Fig. 3) and Cld42+Pr6x7 (x-axis; Fig. 5).

Figure 7. As Fig. 4 but between original Cld42+Pr6x7 (y-axis) and regimes assigned by
precipitation only (x-axis). Regimes with precipitation fractions below 10% have been
combined in the "Others" category.

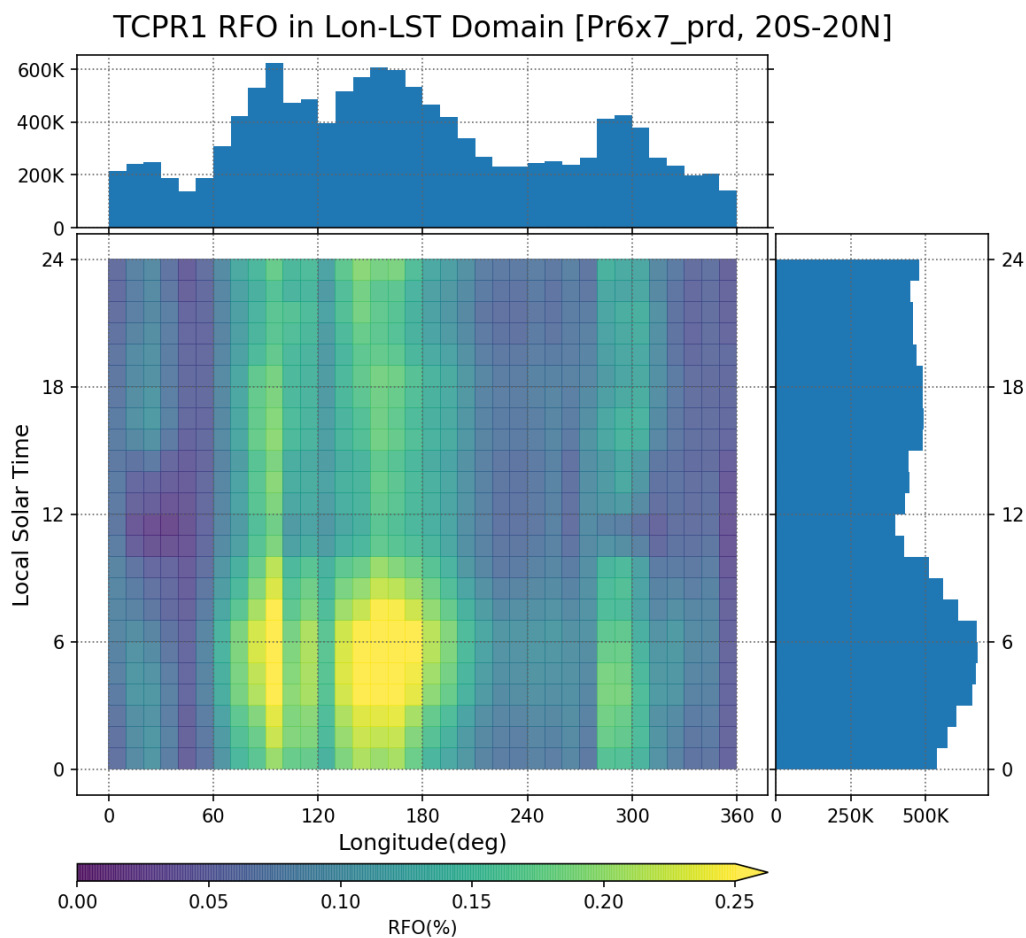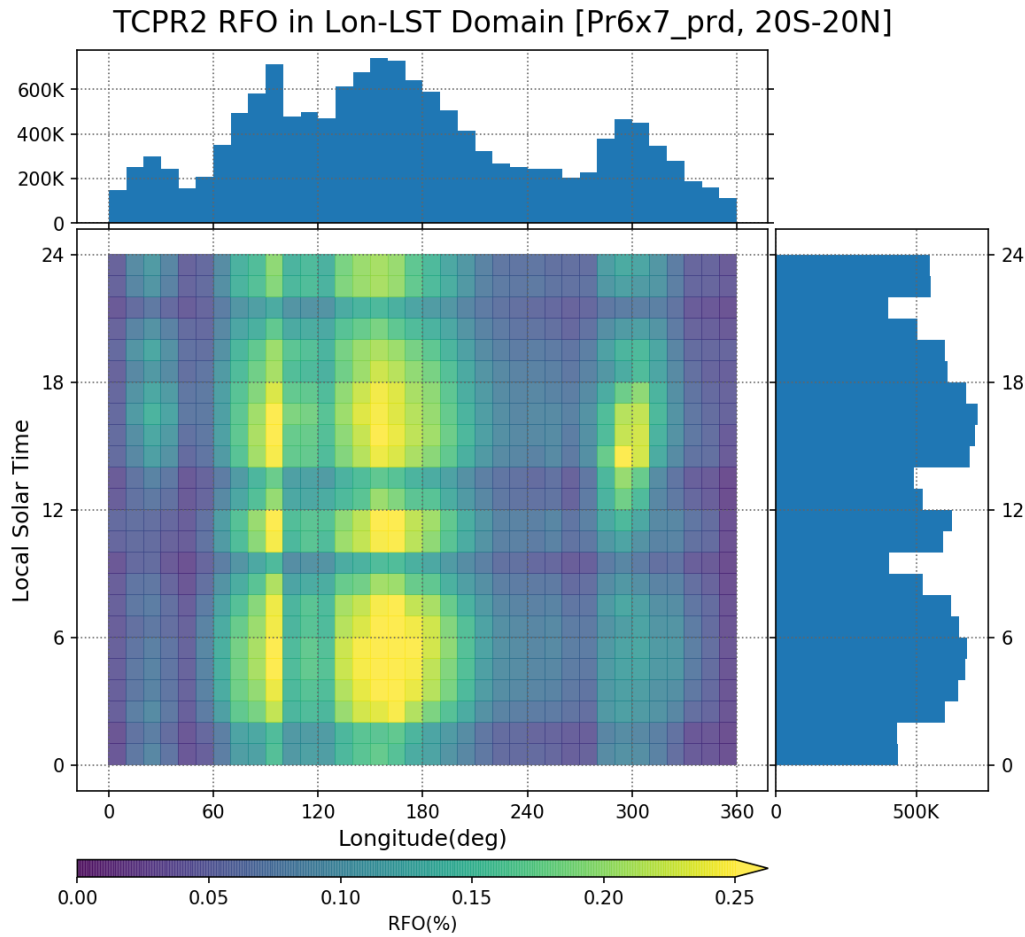Figure 8. Same as Figure 7 but for the set of Cld42+Pr6x1.

754 Figure 9. RFO of TCPR1 of the Cld42+Pr6x7 set predicted by precipitation-only in a
755 longitude (x-axis) and local solar time (LST; y-axis) phase space. Bin resolutions are 10° in
756 longitude, and 1-hour in time. The top and right panels show RFO marginal histograms (sums
757 across rows and columns before normalization) for the same resolution of longitude and LST.

758



TCPR2 RFO in Lon-LST Domain [Pr6x7_prd, 20S-20N]

759

760 Figure 10. Same as Fig. 9, but for TCPR2.

761